

Massimo Alioto

This chapter addresses the challenges and the opportunities to perform computation with nearly-minimum energy consumption through the adoption of logic circuits operating at near-threshold voltages. Simple models are provided to gain an insight into the fundamental design tradeoffs. A wide set of design techniques is presented to preserve the nearly-minimum energy feature in spite of the fundamental challenges in terms of performance, leakage and variations. Emphasis is given on debunking the incorrect assumptions that stem from traditional low-power common wisdom at above-threshold voltages.

In this analysis, the main emphasis is given on the energy consumption, as performance requirements in IoT nodes are easily achievable with near-threshold circuits in most cases, as discussed in Chap. 1 and in the following. Sustained higher levels of performance can always be achieved through architectural techniques (see Chap. 3), whereas occasional performance boosts can be obtained through circuit techniques (see below).

4.1 Preliminary Considerations on Near-Threshold Operation

4.1.1 Transistor Current vs. Supply Voltage and Transregional Model

Voltage scaling is well known to be a very effective knob to reduce the energy per computation at the cost of degraded performance (Burd et al. 2015). The performance degradation at supply voltages V_{DD} lower than the nominal voltage is determined by the reduction in the transistor on-current I_{on} , which in turn depends on the operating region (i.e., the voltage range). The transregional EKV model can be conveniently used to express such dependence in all regions (Enz and Vittoz 2006):

$$I_{on} = I_0 \cdot IC = I_0 \cdot [\ln(e^v + 1)]^2 \quad (4.1)$$

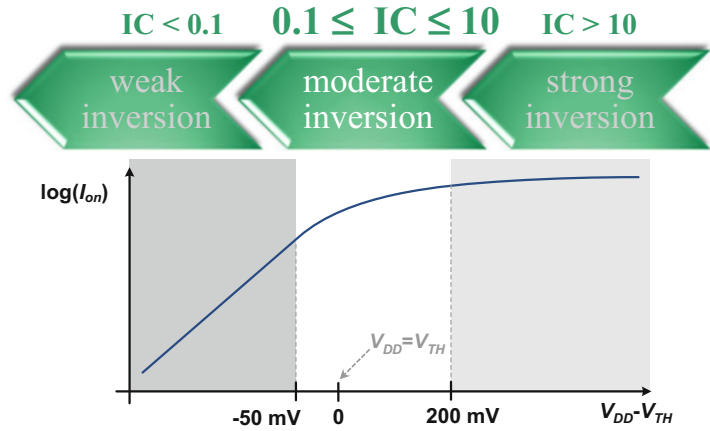
where IC is the inversion coefficient (i.e., normalized current), I_0 is the specific current $2 \cdot n \cdot \mu \cdot C_{OX} \frac{W}{L} (kT/q)^2$, and v is the normalized gate overdrive $v = (V_{DD} - V_{TH})/[2 \cdot n \cdot (kT/q)]$. In the above equations, n is the transistor sub-threshold factor, μ is the carrier mobility, C_{OX} is the MOS capacitance per unit area, W/L is the aspect ratio, V_{TH} is the transistor threshold voltage, and kT/q is the thermal voltage.

In the EKV model in (4.1), a transistor operates in weak inversion when $IC < 0.1$

M. Alioto (✉)

National University of Singapore, Singapore, Singapore
e-mail: malioto@iee.org

Fig. 4.1 Qualitative trend of I_{on} transistor current (log scale) versus the gate overdrive $V_{DD} - V_{TH}$



(i.e., for $v < -1$), which from (4.1) corresponds to voltages below $V_{TH} - 50$ mV for typical n (~ 1.3 – 1.5) and operating temperatures (Sansen 2006). On the other hand, a transistor operates in strong inversion for $IC > 10$ (i.e., for $v > 3.1$), and hence for voltages above $V_{TH} + 200$ mV (Sansen 2006). Near-threshold operation occurs for intermediate voltages, as summarized in Fig. 4.1.

The above traditional EKV model is very useful for quick estimates, but it oversimplifies the I–V characteristics at voltages above V_{TH} . Indeed, eq. (4.1) leads to $I_{on} \approx I_0 \cdot v^2$ in strong inversion, and its quadratic trend is far from the linear trend that is observed in actual nanometer CMOS technologies.¹

Introducing voltage-dependent coefficients in (4.1) solves the issue, but leads to impractically complicated expressions for pencil-and-paper evaluations. To retain its simplicity while employing constant coefficients, (4.1) is here modified according to

$$I_{on} = I_0 \cdot \ln \left(e^{\frac{V_{DD} - V_{TH}}{n \cdot (kT/q)}} + 1 \right) \quad (4.2)$$

which is plotted in Fig. 4.2 along with the actual I–V characteristics for 28-nm NMOS and PMOS transistors. The model is 10% (20%) within

¹ Indeed, sub-100 nm CMOS technologies typically have an I–V characteristics that is proportional to $(V_{DD} - V_{TH})^\alpha$ with $\alpha \approx 1$ (Sakurai and Newton 1990).

circuit simulations on average (in the worst case), hence it is well suited for quick estimates and design purposes.

4.1.2 Transistor Current and Gate Delay in Different Regions

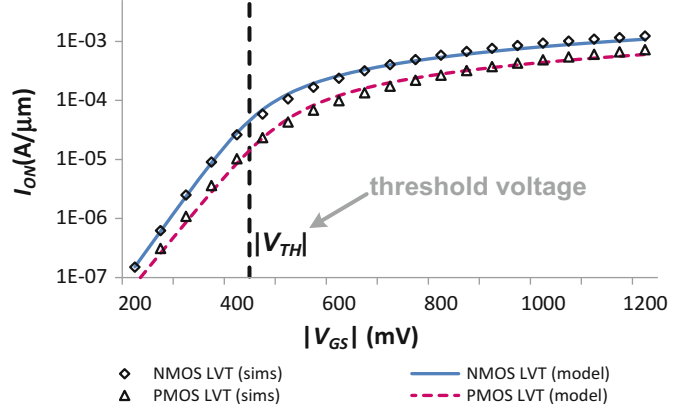
By the definition summarized in Fig. 4.2, sub-threshold voltages correspond to transistor operation in weak inversion, above-threshold are associated with strong inversion, and near-threshold voltages correspond to intermediate voltages between $V_{TH} - 50$ mV and $V_{TH} + 200$ mV. For typical standard threshold voltages,² near-threshold voltages are in the range of 400–600 mV, approximately.

At above-threshold voltages such that $e^{\frac{V_{DD} - V_{TH}}{n \cdot (kT/q)}} \gg 1$, Eq. (4.2) is approximately a linear function of $V_{DD} - V_{TH}$ as expected

² Operation at near-threshold voltages tends to increase V_{TH} compared to the value at nominal voltage, due to DIBL (see Sect. 5.2.2). For standard V_{TH} of 350–380 mV at nominal voltage, it is common to have V_{TH} in the order of 400–450 mV when operating at near-threshold voltages (see, e.g., Fig. 4.7). Observe that the “standard V_{TH} ” nomenclature might be attributed to different threshold voltages in some processes.

Fig. 4.2 Plot of I_{on} transistor current (log scale) in (4.2) versus the magnitude of the gate-source voltage V_{GS} (in CMOS logic gates, $V_{GS} = V_{DD}$)

VOLTAGE RANGE	sub threshold	near threshold	above threshold
TRANSISTOR REGION	weak inversion	moderate inversion	strong inversion
EKV MODEL $IC=I/I_0$ RANGE	$IC < 0.1$	$0.1 < IC < 10$	$IC > 10$
V_{DD} RANGE	$< V_{TH} - 50$ mV	intermediate	$> V_{TH} + 200$ mV



$$I_{above-threshold} \approx \left(I_0 / n \frac{kT}{q} \right) \cdot (V_{DD} - V_{TH}) \quad (4.3)$$

whereas at sub-threshold voltages it can be approximated as³

$$I_{sub-threshold} \approx I_0 \cdot e^{\frac{V_{DD} - V_{TH}}{n \cdot V_T}} \quad (4.4)$$

which exponentially decreases when lowering the voltage. At near-threshold voltages, Eq. (4.2) can be approximated as

$$I_{near-threshold} \approx \frac{I_0}{2} \cdot \left[1.5 + \left(\frac{V_{DD} - V_{TH}}{n \cdot kT/q} \right)^{1.35} \right] \quad (4.5)$$

which is within 15% of the exact I–V characteristics in Fig. 4.2. From (4.5), the near-threshold I–V characteristics is a power law, and is steeper than in the above-threshold region.

Let us now consider a CMOS logic gate driving a capacitive load C , which includes the capacitive parasitics of the gate itself. As usual

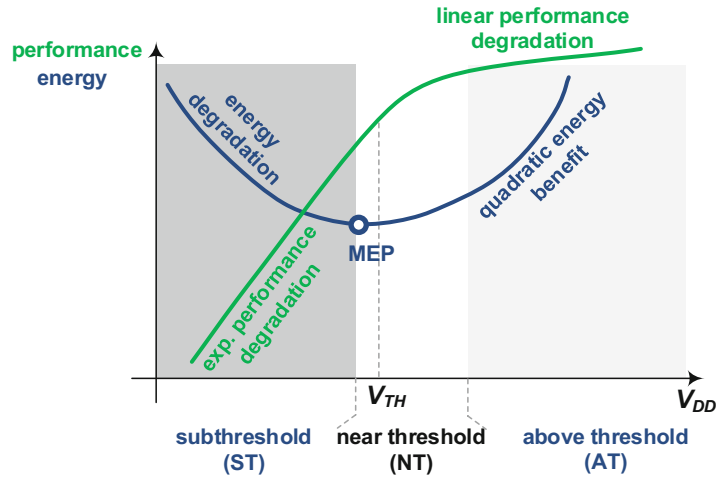
(Weste and Harris 2011), its propagation delay τ_{PD} can be expressed as $(C/I_{on}) \cdot (V_{DD}/2)$:

$$\tau_{PD} = \frac{C}{I_{on}} \cdot \frac{V_{DD}}{2} \approx \frac{C}{2 \cdot I_0} \cdot \frac{V_{DD}}{\ln \left(e^{\frac{V_{DD} - V_{TH}}{n \cdot (kT/q)}} + 1 \right)}. \quad (4.6)$$

As shown in Fig. 4.3, from (4.3) and (4.6), voltage downscaling leads to an approximately linear delay (i.e., performance) degradation, when operating above threshold. As discussed in Sect. 4.3, the energy is typically dominated by the dynamic contribution, hence a quadratic energy saving is observed above threshold. On the other hand, an exponential increase in the gate delay is observed in the sub-threshold region. Also, due to the heavier leakage contribution at low voltages (Sect. 4.3), the energy reaches a minimum energy point (MEP), and it tends to increase again when further lowering V_{DD} . Hence, near-threshold voltages are an ideal compromise between energy and performance in energy-centric VLSI designs. Indeed, the near-threshold gate delay is still reasonably small, and energy is close to its minimum value across all voltages. This motivates this chapter, and the adoption of near-threshold circuits for VLSI processing in the IoT domain.

³ Indeed, $\ln(e^x + 1) \approx e^x$ for $x < 0$ (i.e., for $V_{DD} < V_{TH}$) in (5.2).

Fig. 4.3 Qualitative trend of performance (gate delay) and energy per operation versus supply voltage V_{DD}



4.2 Near-Threshold Transistor and Circuit Properties

In this section, properties of transistors at near threshold are discussed to provide general circuit design guidelines. Preliminary considerations on voltage scaling and threshold voltage dependence on sizing are respectively provided in Sects. 4.2.1 and 4.2.2. The impact of transistor stacking and PMOS/NMOS imbalance are discussed in Sect. 4.2.3 to guide the topology selection during circuit design. As second and equally fundamental aspect of circuit design, transistor strength adjustment is discussed in Sect. 4.2.4.

4.2.1 Impact of Aggressive Voltage Scaling on Transistor Current and Delay

The considerations on the delay degradation under voltage scaling in the previous section were based on the assumption that the gate load C is independent of V_{DD} . Observe that the load C comprises wire parasitics and transistor gate capacitances. The above assumption certainly holds in wire-dominated loads (as wire parasitics are voltage-independent), whereas it is somewhat pessimistic in gate-dominated loads. Indeed, as

shown by Fig. 4.4, the transistor gate capacitance tends to moderately decrease at voltages close to or below V_{TH} , and hence makes the delay degradation more graceful than discussed above, although to a minor extent.

According to the above observation, the above qualitative considerations on the delay at near- and sub-threshold voltages fully apply to any practical design. As an example, Fig. 4.5 shows the trend of the fan-out-of-4 delay $FO4$ (i.e., the delay of an inverter gate driving four equal inverters). This metric is widely used at process level to characterize the speed of the technology, at circuit level to abstract the circuit design from the process details, and at architectural level since the clock cycle normalized to $FO4$ is typically a constant that is defined by the architecture (Harris). In short, $FO4$ characterizes the system performance versus voltage for a given architecture. From Fig. 4.5, operation in the middle of the near-threshold region degrades the performance by approximately a factor of 10, compared to operation at nominal voltage. This is generally true regardless of the adopted technology (Dreslinski et al. 2010).

A very distinctive property of near-threshold operation is the stronger delay sensitivity to a given absolute change in the gate overdrive (i.e., both V_{DD} and V_{TH}), compared to above-threshold designs. This is partially explained by the steeper I - V characteristics (4.5) compared to (4.3) (the exponent of v is respectively 1.35 and 1). But the

Fig. 4.4 Gate capacitance normalized to value at nominal voltage versus supply voltage V_{DD} (28 nm, LVT and RVT transistors)

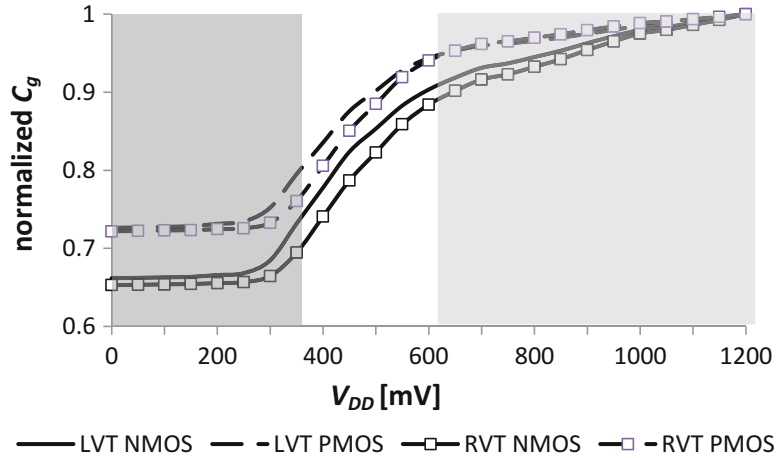
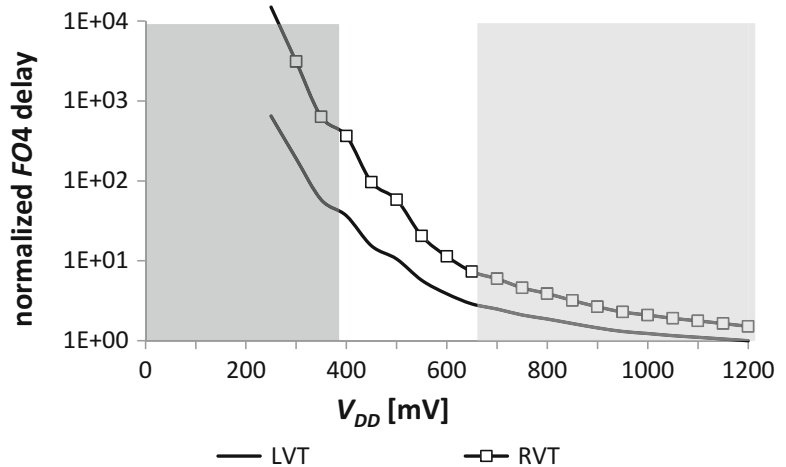


Fig. 4.5 Fanout-of-4 (FO4) delay normalized to value at nominal voltage versus supply voltage V_{DD} (28 nm, LVT and RVT transistors)



main reason is due to the very large sensitivity of $v = (V_{DD} - V_{TH})/[2 \cdot n \cdot (kT/q)]$ to a given change in V_{DD} , as V_{DD} is much closer to V_{TH} compared to above-threshold voltages. The relative I_{on} improvement due to a 100-mV supply voltage increase (i.e., boosting) for a 28-nm technology is shown in Table 4.1. As expected, the impact of voltage boosting at near-threshold voltages is substantially larger than above threshold, with improvements in I_{on} in the range of 2-4X. This unique feature permits to have significant speed adjustment capability with very limited amount of boosting, which needs to be thoroughly exploited in near-threshold-designs.

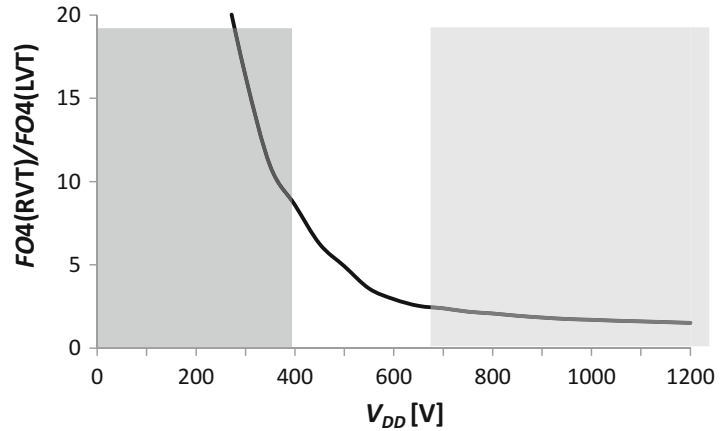
The above considerations equally apply to the threshold voltage, as I_{on} is a direct function of the

Table 4.1 I_{on} Improvement due to supply voltage boosting by 100 mV

V_{DD}	I_{on} improvement
400 mV	4.05X
500 mV	2.24X
600 mV	1.7X
800 mV	1.31X
1 V	1.17X

gate overdrive $V_{DD} - V_{TH}$. For example, the I_{on} and speed sensitivity to a 100-mV V_{DD} shift in Table 4.1 hold for the same change in V_{TH} (although with negative sign). In other words, increasing V_{TH} by 100 mV (i.e., the typical difference between a low and regular V_{TH}) at near-threshold supply voltages leads to a 2-4X

Fig. 4.6 Ratio between $FO4$ of RVT and LVT transistors versus supply voltage V_{DD}



reduction in speed. This is shown in Fig. 4.6, which plots the inverter delay ratio under regular- V_{TH} (RVT) and low- V_{TH} (LVT). At nominal voltage, the different V_{TH} has a moderate impact on the performance, whereas such difference is much more pronounced at near-threshold voltages.

In summary, the high sensitivity of performance to V_{DD} and V_{TH} makes them very powerful knobs at near-threshold voltages, although the (same) sensitivity to their variations poses a challenge at the same time, as will be discussed in the following sections.

4.2.2 Impact of DIBL and Sizing on Threshold Voltage

In the previous subsection, V_{TH} was implicitly considered constant. In view of the large sensitivity of I_{on} to V_{TH} , the dependence of V_{TH} on transistor voltages and sizing needs to be explicitly considered at near-threshold voltages.

Regarding the dependence of the transistor voltages, V_{TH} tends to be quite sensitive to the drain-source voltage due to the Drain Induced Barrier Lowering (DIBL) effect (Tsividis 1999). Due to the DIBL effect, V_{TH} increases in an approximately linear fashion when the magnitude of the drain-source voltage is reduced. Due to the body effect, V_{TH} decreases (increases) under Forward FBB (Reverse, RBB) Body

Biasing, i.e. for positive⁴ (negative) body bias voltages V_{BB} (Tsividis 1999). The approximately linear dependence in both effects is captured by the following equation

$$V_{TH} = V_{TH0} - \lambda_{DIBL}V_{DS} - \lambda_{BB}V_{BB} \quad (4.7)$$

where V_{TH0} is the threshold voltage extrapolated for very low V_{DD} and $V_{BB} = 0$ V, λ_{DIBL} is the DIBL coefficient and λ_{BB} is the body effect coefficient. The DIBL coefficient is in the order of 0.1 V/V or larger for technologies suited for IoT, and hence denotes a pronounced dependence of the threshold voltage on the drain-source voltage. As an example, Fig. 4.7 shows the change in V_{TH} versus the drain-source voltage (i.e., V_{DD}) in 28-nm transistors. When V_{DD} is reduced down to near-threshold voltages, V_{TH} typically increases by around 100 mV compared to operation at nominal voltage. This needs to be explicitly taken into account when choosing the type of threshold voltage at design time.

On the other hand, the threshold voltage dependence on the body voltage is well known to be rather weak in advanced technologies, although it is appreciable in 90-nm generations or older. Considering the strong sensitivity of performance and leakage on V_{TH} , body biasing

⁴These considerations hold for NMOS transistors. For PMOS transistors, change the sign in all voltages. Regarding the body effect, FBB (RBB) refers to body voltages V_{BB} below (above) V_{DD} .

Fig. 4.7 Threshold voltage deviation vs. drain-source voltage due to DIBL

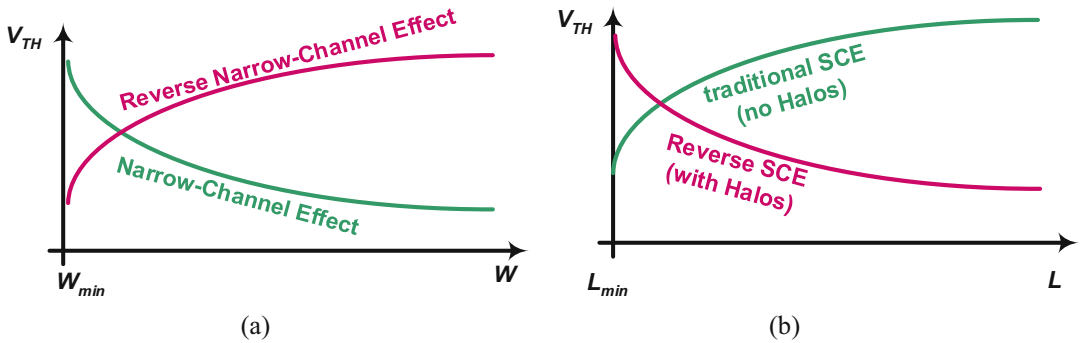
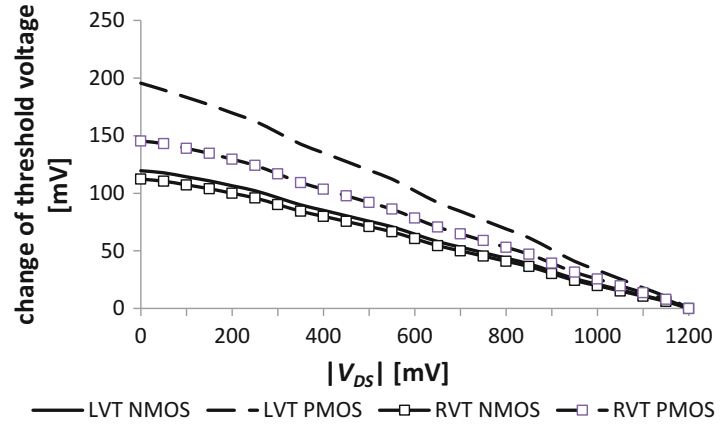


Fig. 4.8 Qualitative trend of threshold voltage vs. (a) transistor channel width, (b) transistor channel length (due to Short Channel Effect)

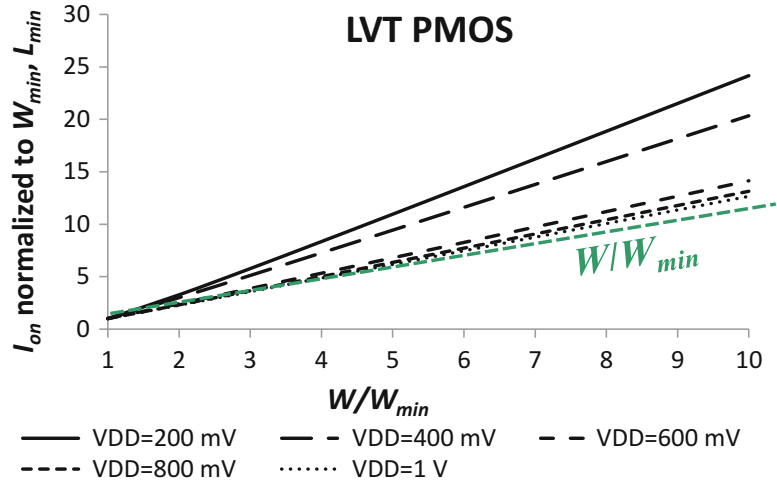
is still a viable option in near-threshold circuits in 90-nm technology bulk generations, or in more recent generations in FDSOI CMOS technology.

The transistor threshold voltage also depends on the size, especially when the latter is close to the minimum allowed by the process. The qualitative dependence on the channel width W (length L) is qualitatively depicted in Fig. 4.8a, b. From Fig. 4.8a, the reduction of W leads to a decrease (increase) in V_{TH} due to the Reverse (traditional) Narrow-Channel Effect RNCE (NCE) (Tsividis 1999). The dominance of one of the two effects mainly depends on the transistor isolation technology (e.g., Shallow Trench Isolation vs LOCOS), device structure (bulk, FinFET, FDSOI) and several parameters. On the other hand, from Fig. 4.8b, the reduction

of L leads to an increase (decrease) in V_{TH} due to the Reverse (traditional) Short-Channel Effect RSCE (SCE) (Tsividis 1999). The dominance of one of the two former mainly depends on whether the transistor body is lightly doped or includes halos to counteract short-channel effects (Tsividis 1999).

From the above considerations, transistor sizing can affect the performance in ways that are more complicated than the usual linear dependence of I_{on} on W/L , due to the additional (strong) dependence of V_{TH} on size at near-threshold voltages. As an example, Fig. 4.9 shows that the I_{on} current trend versus W deviates from the traditional linear dependence of $I_{on} \propto W$. For the specific considered technology, the current is increasing faster than

Fig. 4.9 I_{on} normalized to current in minimum-sized transistor vs. channel width (normalized to minimum width W_{min})



$I_{on} \propto W$, denoting that the NCE effect dominates (i.e., wider channels lead to large current than expected due to the simultaneous reduction in V_{TH}). This effect is clearly more pronounced at low voltages due to the stronger dependence of I_{on} on V_{TH} , whereas it is negligible at nominal voltage.

Other technologies might have opposite behavior due to dominant RNCE (i.e., I_{on} increases slower than W , due to the progressive increase in V_{TH} due to the increase in V_{TH}). On the other hand, Fig. 4.10 shows that I_{on} decreases faster than $1/L$ at near-threshold voltages, due to the dominance of SCE. Again, this dependence is 1.5–3X stronger than at nominal voltage due to the stronger dependence of I_{on} on V_{TH} at low voltages. Other technologies might have different behavior, due to the dominance of RSCE.

4.2.3 PMOS/NMOS Strength Ratio, Stacking and Wire Delay

Another important effect observed at near-threshold voltages is the deviation of the ratio of the PMOS and NMOS strength (i.e., I_{on}) at iso-size, compared to nominal voltage. This is due to the different dependence of PMOS and NMOS I_{on} across different voltages. Indeed, from (4.3)–(4.6) the transistor strength has a

mild dependence on V_{TH} and is mostly defined by the carrier mobility at nominal voltage. Hence, differences in V_{TH} between PMOS and NMOS do not significantly impact the strength. On the other hand, the strength has a strong dependence on V_{TH} at near-threshold (and lower) voltages, hence even moderate differences in V_{TH} between PMOS and NMOS substantially alter their strength ratio. The latter can be smaller or larger than the value at nominal voltage, depending on the V_{TH} differences between PMOS and NMOS (including DIBL), and hence the specific technology. Figure 4.11 shows the trend of the PMOS/NMOS strength ratio in a specific 28-nm technology, which at near-threshold voltages can be reduced by up to 2.5X compared to nominal voltages. At lower voltages, the impact is even larger, due to the exponential dependence of I_{on} on V_{TH} in (4.4). This deviation of the PMOS/NMOS strength ratio clearly threatens the noise margin of CMOS logic gates, thus degrading robustness and exposing logic gates to malfunctions due to variations. This also emphasizes the imbalance between the rise and fall delay, thus degrading performance.

Analogously, the strength of stacked transistors (i.e., connected in series) can heavily deviate from the strength of a single transistor, compared to operation at nominal voltage.

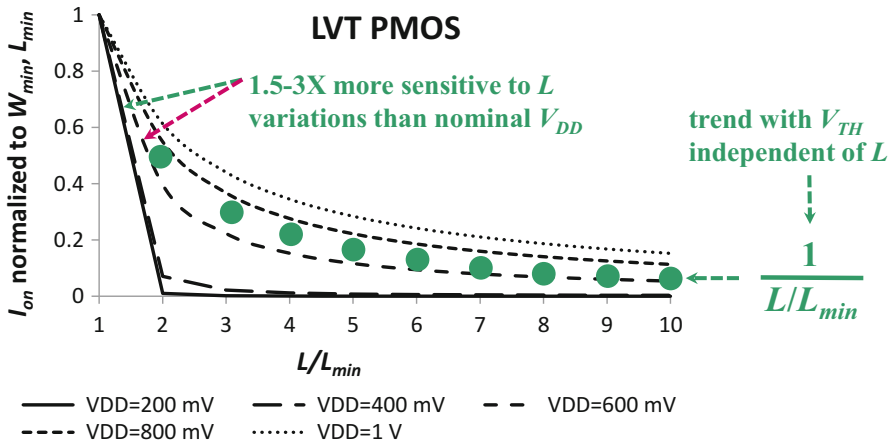
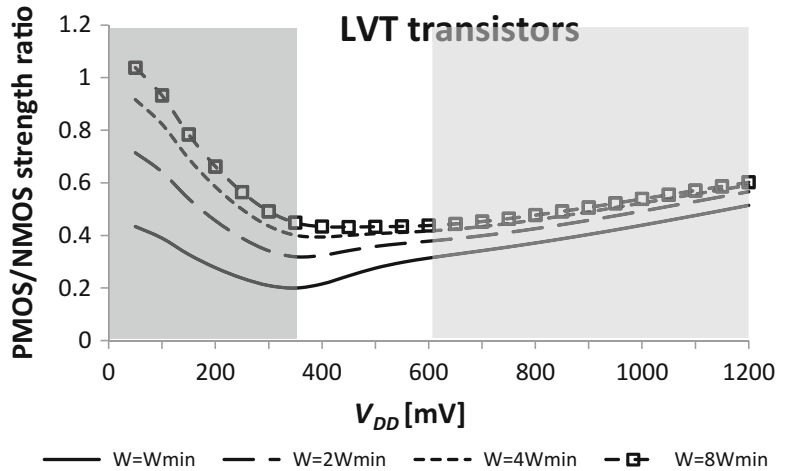


Fig. 4.10 I_{on} normalized to current in minimum-sized transistor vs. channel length (normalized to minimum length L_{min})

Fig. 4.11 Ratio between the strength of PMOS and NMOS versus supply voltage V_{DD} (LVT transistors)



This can be shown by the I_{on} stacking factor X_{on} , defined as the factor by which the I_{on} current is reduced due to the transistor stacking, compared to a single transistor (assuming all transistors have the same size as the single one). The trend in Fig. 4.12 shows that the stacking factor tends to peak around near-threshold voltages, and the phenomenon is more evident under a larger number of stacked transistors. At lower (sub-threshold) voltages, the stacking factor goes back to smaller and threshold-voltage independent value (Alioto 2012).

The stacking factor peaking at near-threshold voltages can be observed in any CMOS technology, as the presence of stacked transistors reduces the drain-source voltage of each stacked transistor, and hence leads to a further increase in V_{TH} (and decrease in the strength) due to DIBL, compared to a single transistor. This explains why the near-threshold current delivered by four stacked transistors is up to 7X lower than a single transistor at iso-size, although this factor is about half of it at nominal voltage. Due to the same reason, the degradation

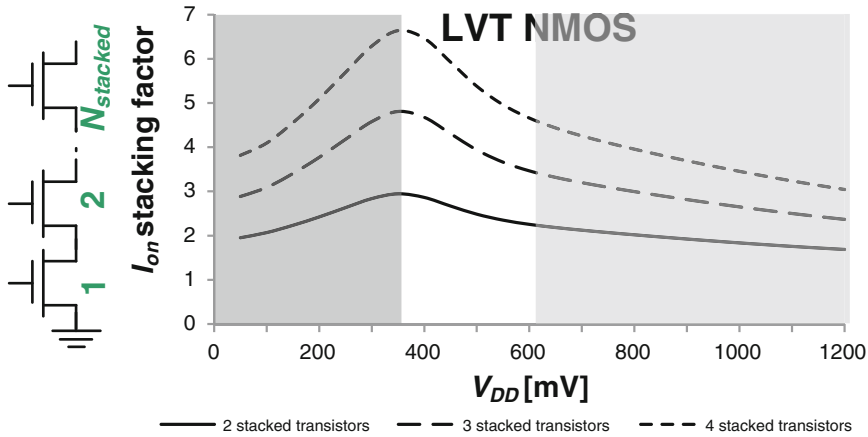


Fig. 4.12 Ratio between the strength of PMOS and NMOS versus supply voltage V_{DD} (LVT transistors)

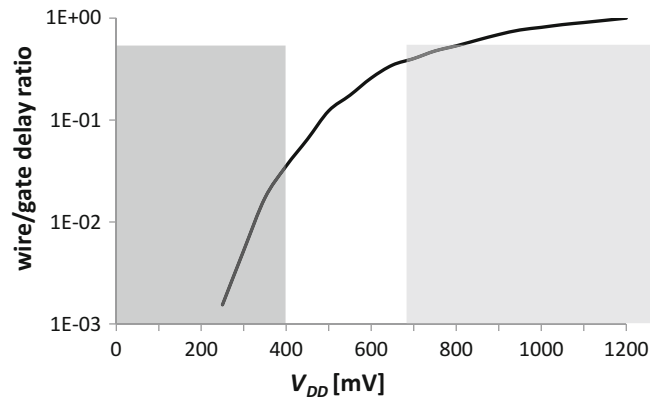
for two and three stacked transistors is much less pronounced, and is acceptable from a performance point of view. Hence, as general circuit design guideline, the maximum number of stacked transistors in near-threshold designs needs to be lower (e.g., 3) than at nominal voltage (typically up to four).

Finally, another fundamental difference encountered in near-threshold designs is the deviation of the ratio between the gate and wire delay, compared to nominal voltage. Indeed, at lower voltages, the gate delay increases as in (4.6), whereas the wire delay remains constant. As an example, Fig. 4.13 considers a wire whose delay matches the delay of a single gate designed for high performance (i.e., its delay is about one $FO4$ (Sutherland et al. 1999)) at nominal V_{DD} . This corresponds to a global wire with a length in the order of a millimeter. From this figure, the wire delay at near-threshold voltages represents only a small fraction (in the order of 10X smaller) of the gate delay. This means that above-threshold designs and architectures that aims at mitigating the impact of wire delay are definitely overdesigned and performance/energy sub-optimal at near-threshold voltages. Hence, near-threshold circuits and architectures need to be very different from traditional above-threshold solutions, due to drastically smaller

impact of wire delay. This brings back circuit and architectural solutions that were abandoned in the late 90s, due to the then incumbent impact of wires on the system performance.

In summary, the large performance sensitivity on V_{DD} and V_{TH} represents a very interesting opportunity in near-threshold circuits, but also poses various challenges. Among those, it is not possible to maintain a fixed delay ratio between cells with different amount of stacking and threshold voltage, when scaling the voltage (even without variations). This, in addition to the substantial performance degradation due to stacked transistors, suggests that the maximum fan-in of near-threshold CMOS standard cell should be three. Within the same cell, it is not possible to maintain a stable PMOS/NMOS strength ratio across different voltages. For the same reasons, ratioed and dynamic logic styles are unfeasible at near-threshold voltages (not to mention the larger impact of leakage and variations, as discussed in Sects. 4.3 and 4.6). Similarly, topologies that are inherently based on current contention and positive feedback need to be definitely avoided (e.g., cross-coupled non-clocked inverters in flip-flops). Unfortunately, this cannot be avoided in SRAM bitcells and register files for reasons due to density, and other sophisticated techniques need to be deployed (see Chap. 5).

Fig. 4.13 Ratio of wire and gate delay normalized to value at nominal voltage versus supply voltage V_{DD} (28 nm, LVT transistors)



4.2.4 Knobs to Adjust Transistor Strength

From the previous subsection, the transistor strength can be adjusted with the following knobs:

- transistor size
- body biasing
- V_{TH} selection
- V_{DD} tuning and fine-grain boosting.

From the previous subsection, transistor sizing is relatively effective, and can be more or less effective than at nominal voltage, depending on the dominance of RNCE over NCE, and SCE over RSCE. Body biasing can significantly alter the transistor strength only in old technologies (e.g., 90 nm), or in recent FDSOI technologies, with a typical 30% range of adjustment at near-threshold voltages.

In view of the strong dependence of I_{on} on the gate overdrive discussed in Sect. 4.1, the transistor strength can be substantially modified through the proper selection of the threshold voltage, and the fine-grain boosting of V_{DD} to selectively increase I_{on} where required. Regarding the V_{TH} selection, a 2–4X I_{on} (and delay) change was previously shown to be feasible when changing V_{TH} from one type (e.g., RVT) to the next available one (e.g., LVT). However, V_{TH} selection at near-threshold voltages poses various additional challenges, compared to

operation at nominal voltage. Indeed, the sensitivity of I_{on} to V_{TH} translates into a strong sensitivity to its process variations. Also, the delay ratio of an RVT and LVT logic gate (see Fig. 4.14. for an inverter gate) strongly depends on the supply voltage. In other words, mixing standard cells with different V_{TH} poses the problem of having different delay scaling in different portions of the system. In turn, this makes timing closure certainly more difficult and might reduce the energy benefit of dynamic voltage scaling, as the critical path(s) depends on the voltage.

Let us now consider fine-grain voltage boosting, which consists in selectively overdriving appropriate transistors with a voltage above V_{DD} . As shown in the illustrative example in Fig. 4.15, this might be the case of a single large transistor M1 (e.g., sleep transistor, large buffer) that drives a sub-circuit containing several smaller transistors. Let us assume that the gate of M1 is overdriven at $V_{DD} + \Delta V_{DD}$ as opposed to all other transistors and logic gates, which are powered at V_{DD} . Due to the strong (super-linear) I_{on} increase in M1 due to the gate voltage boosting by ΔV_{DD} , the transistor can be significantly undersized while maintaining the same strength as the transistor that is driven by V_{DD} . In view of the strong dependence of I_{on} on the gate voltage in M1 at near-threshold voltages, a small amount of boosting ΔV_{DD} permits to substantially reduce the area occupied by M1. This is shown in the example in Fig. 4.15 in 28 nm, where the area of M1 can be reduced by

Fig. 4.14 Ratio of I_{on} of LVT and RVT transistors normalized to value at nominal voltage versus supply voltage V_{DD}

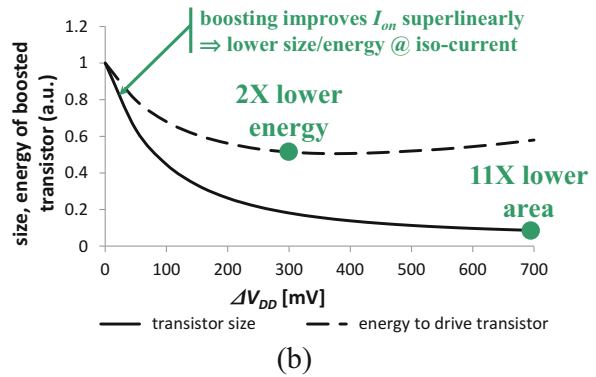
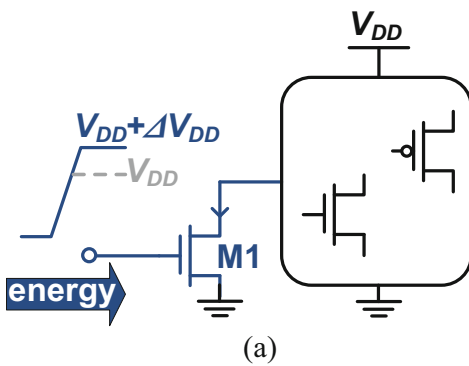
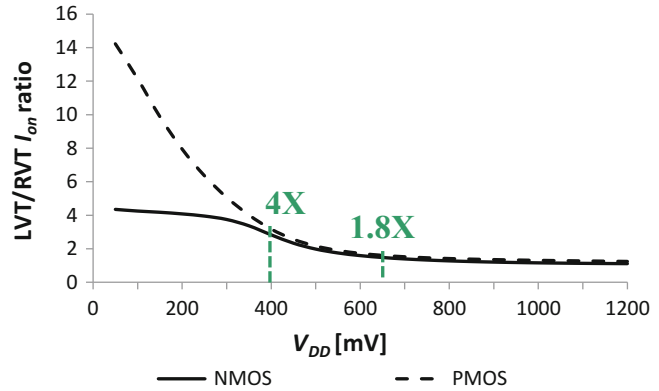


Fig. 4.15 (a) In-principle circuit with large transistor whose gate voltage is boosted by ΔV_{DD} , (b) area and energy improvement vs. ΔV_{DD}

up to an order of magnitude while maintaining the same strength, through an amount of boosting in the order of a few hundreds of mVs. Similarly, such selective boosting permits to super-linearly reduce the leakage current of M1. At the same time, the gate capacitance $C_{g,M1}$ of M1 is reduced super-linearly as well, whereas the supply voltage is increased by the very limited amount ΔV_{DD} . This means that the dynamic energy $C_{g,M1} \cdot (V_{DD} + \Delta V_{DD})^2$ to switch M1 ON is reduced overall. In the example in Fig. 4.15, a 2X energy reduction can be achieved through selective boosting of M1, when adopting an optimal ΔV_{DD} of 300 mV (this voltage depends on the specific technology). Very similar energy saving is observed for ΔV_{DD} in the order of 100–200 mV. On the other hand, larger amount of boosting slightly increases the energy consumption to turn on M1, since the transistor starts

operating above threshold (i.e., I_{on} becomes less sensitive to ΔV_{DD}), and the energy cost $C_{g,M1} \cdot (V_{DD} + \Delta V_{DD})^2$ of boosting increases substantially due to the quadratic dependence.

From the above considerations, near-threshold circuits can be made more energy- area-efficient by selectively boosting portions of the circuit that contain large (and hence energy- and area-hungry) transistors. As opposed to traditional multi- V_{DD} approaches that are applied at the module level, in this case the supply is boosted with fine granularity (i.e., down to the single transistor). Such fine-grain voltage boosting also offers the opportunity to equalize imbalanced logic across pipestages. As fundamental challenge, fine-grain boosting entails significant area overhead, due to the additional level shifters to drive boosted-voltage domains, and to the slight additional cost of distributing multiple voltages at

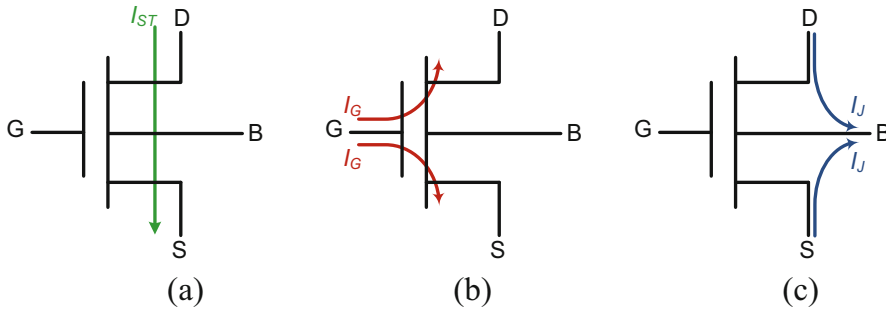


Fig. 4.16 Transistor leakage contributions: (a) sub-threshold, (b) gate, (c) substrate

the physical design level (Flynn et al. 2007). In other words, near-threshold circuits should certainly take advantage of fine-grain voltage boosting, but innovative techniques are needed to minimize the unavoidable overhead. Some recently proposed ideas to address this challenge will be presented in Sect. 4.1.

4.3 Energy Trends

In this section, the impact of voltage scaling on the energy is reviewed, providing models and design guidelines for minimum-energy operation.

4.3.1 Transistor Leakage Current at Near-Threshold Voltages

The MOS transistor leakage contributions are summarized in Fig. 4.16. The dominant contribution is due to the sub-threshold leakage in (4.4), which flows between drain and source and is due to the diffusion of minority carriers between the two terminals (Tsividis 1999). The gate leakage flows from gate to source/drain or vice versa, depending on the applied voltages, and tends to be exponentially smaller than the sub-threshold contribution when lowering V_{DD} (Narendra and Chandrakasan 2006). Similar considerations hold for the substrate leakage, which is mostly due to the Band-to-Band Tunneling (BTBT), and the inverse saturation current of the source-bulk and drain-bulk pn junctions (Narendra and Chandrakasan 2006). Hence, the transistor leakage current at near-

threshold voltages is well approximated by (4.4), where the gate-source voltage (assigned to V_{DD} in (4.4)) has to be set to zero. By substituting the dependence of V_{TH} in (4.7), the near-threshold leakage current of an NMOS transistor immediately results to

$$I_{lkg} = I_0 \cdot e^{-\frac{V_{TH0} - \lambda_{BB} V_{BB}}{nKT/q}} \cdot e^{-\frac{\lambda_{DIBL} V_{DD}}{nKT/q}} \quad (4.8)$$

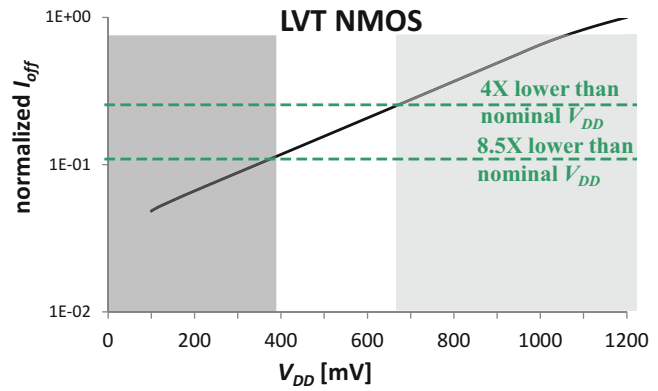
where the first exponential term is set by the threshold voltage (including body biasing, when applied), and the second expresses the DIBL effect and hence the leakage dependence on V_{DD} . The latter dependence is exponential, and typically operation at near-threshold voltages reduces the leakage current by about an order of magnitude for typical DIBL coefficients in the order of 0.1 V/V, compared to nominal voltage. The consistent exponential trend across voltages in (4.8) is shown in Fig. 4.17 for a 28-nm technology, along with a leakage current reduction at near-threshold voltages by 4–8.5X.

4.3.2 Energy Consumption of Digital Systems at Near-Threshold Voltages

The total energy per operation⁵ in a near-threshold VLSI digital system is essentially equal to the sum of the dynamic and the leakage energy. Indeed, the short-circuit energy

⁵ An operation is here defined as the basic task that the considered system is executing, e.g., an instruction in a CPU or GPU, a new output sample in a DSP, an arithmetic operation in an Arithmetic Logic Unit.

Fig. 4.17 Leakage current I_{off} of LVT transistor (normalized to value at nominal voltage) versus supply voltage V_{DD}



contribution (Weste and Harris 2011) is negligible at near-threshold voltages, as opposed to operation at nominal voltage. This is because the transistors for input voltages around $V_{DD}/2$ is a small sub-threshold current, which also rapidly vanishes when the input voltage deviates from $V_{DD}/2$ to settle to its stable value (Alioto 2012) (due to the exponential I–V characteristics in the sub-threshold region).

The dynamic energy per operation is given by

$$E_{dyn} = \alpha_{SW} \cdot C \cdot V_{DD}^2 \cdot CPO$$

✓ ✓ ✓ affected by (micro) architecture
✓ ✓ ✓ affected by circuit design
✓ ✓ affected by technology

$$E_{lkg} = V_{DD} \cdot I_{off} \cdot T_{CK} \cdot CPO$$

where $\alpha_{SW} \cdot C \cdot V_{DD}^2$ is the energy per cycle, being C the total capacitance within the circuit, α_{SW} is the activity factor (Weste and Harris 2011) (i.e., the fraction of C that is switched in a cycle, on average). In (4.9), it was considered that an operation in general takes an average number of cycles CPO (Cycles per Operation), which depends on the specific (micro)architecture, and the dataset to a minor extent (e.g., in microprocessors).

The leakage energy per operation can be expressed as the product of the average leakage power $V_{DD} \cdot I_{off}$ (being I_{off} the average leakage current), the clock cycle T_{CK} and CPO (Alioto 2012). T_{CK} can be expressed as $FO4 \cdot LD_{eff}$, where $LD_{eff} = T_{CK}/FO4$ is the number of the

number of $FO4$ delays (i.e., cascaded inverters with fan-out of 4) that can fit the cycle time. Hence, LD_{eff} represents the effective logic depth per pipestage, which is a constant defined by the (micro)architecture.⁶ Hence the leakage energy per operation can be written as

$$E_{lkg} = V_{DD} \cdot I_{off} \cdot T_{CK} \cdot CPO$$

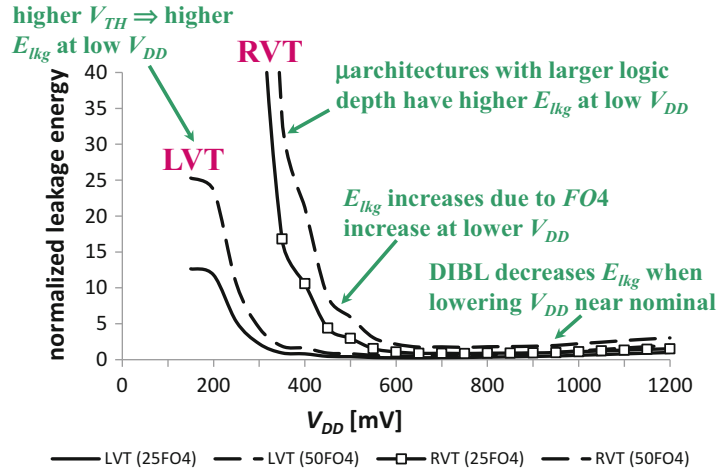
✓ ✓ ✓ affected by (micro)architecture
✓ ✓ ✓ affected by circuit design
✓ ✓ affected by technology

(4.10)

In (4.10), the only parameters that depend on V_{DD} are V_{DD} itself, I_{off} and $FO4$. When down-scaling the voltage, the first term decreases linearly and I_{off} decreases exponentially due to DIBL, although not very rapidly since V_{DD} is multiplied by $\lambda_{DIBL} \ll 1$ in (4.8). On the other hand, $FO4$ rapidly increases as in (4.6) when V_{DD} is reduced down to near-threshold voltages and below. The overall effect of the three factors leads to an increase in E_{lkg} at near- and sub-threshold voltages when V_{DD} is reduced, as opposed to the dynamic energy. The leakage energy tends to increase very rapidly when decreasing V_{DD} down to the transistor threshold

⁶ $T_{CK}/FO4$ is essentially constant in gate-delay dominated critical paths when varying V_{DD} , as all gate delays generally scale like $FO4$ (Harris et al. n.d.; Weste and Harris 2011). At near-threshold voltages, this assumption is generally correct, as the wire delay is typically much smaller (see Sect. 5.2.3).

Fig. 4.18 Leakage energy vs. supply voltage V_{DD} for different logic depths LD_{eff} equal to $25FO4$ and $50FO4$ (28 nm, LVT and RVT transistors)



voltage, due to the resulting rapid increase in the gate delay in (4.6). This is shown in Fig. 4.18, where the leakage energy of the reference digital circuit in Sect. 4.4 is plotted versus V_{DD} under RVT and LVT transistor flavor. As expected, the leakage energy under RVT flavor rapidly increases at higher voltages compared to LVT, due to the higher threshold voltage. This figure also shows that E_{lkg} tends to shoot up at larger voltages, under microarchitectures with larger logic depth (e.g., $LD_{eff} = 50$ instead of 25). This is because such microarchitectures suffer from larger leakage energy from (4.10), and hence the rapid increase can be observed at larger voltages. On a side note, Fig. 4.18 also shows that E_{lkg} has an opposite behavior at above-threshold voltages (i.e., it decreases when decreasing V_{DD}), due to the dominance of the exponential effect of DIBL over the linear $FO4$ increase.

From (4.9)–(4.10), the total energy per operation E_{TOT} of a given VLSI system or sub-system results to

$$E_{TOT} = E_{dyn} + E_{lkg} = E_{cycle} \cdot CPO \quad (4.11)$$

where the energy per cycle $E_{cycle} = E_{TOT}/CPO$ is defined as

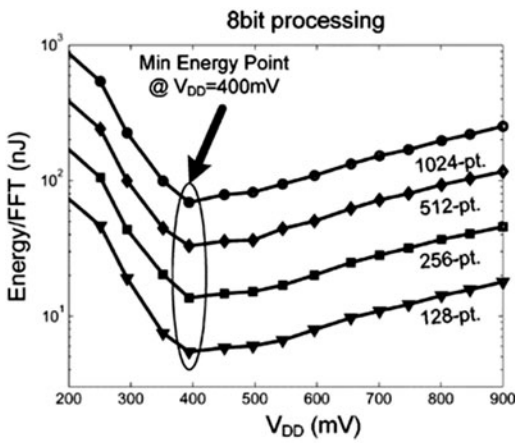
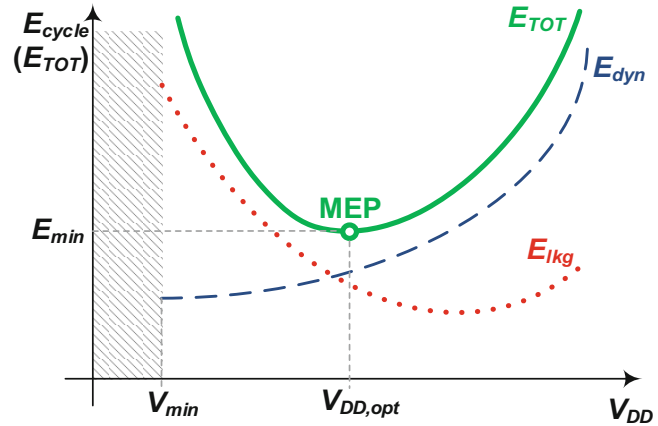
$$E_{cycle} = \alpha_{SW} \cdot C \cdot V_{DD}^2 + V_{DD} \cdot I_{off} \cdot FO4 \cdot LD_{eff} \quad (4.12)$$

The qualitative trend of (4.11)–(4.12) versus V_{DD} in Fig. 4.19 shows that the voltage down-scaling reduces the dynamic energy, but increases the leakage energy. Hence, a minimum-energy point (MEP) is observed at a voltage $V_{DD,opt}$ that optimally balances the dynamic and leakage energy, thus leading to the minimum energy⁷ E_{min} . The MEP voltage $V_{DD,opt}$ typically lies in the sub-threshold or near-threshold region (Hanson et al. 2006a; Hanson et al. 2006b), as discussed in the next section. Due to the flatness of the MEP, near-threshold operation permit true- or nearly-minimum energy operation, as fundamental design target of this chapter.

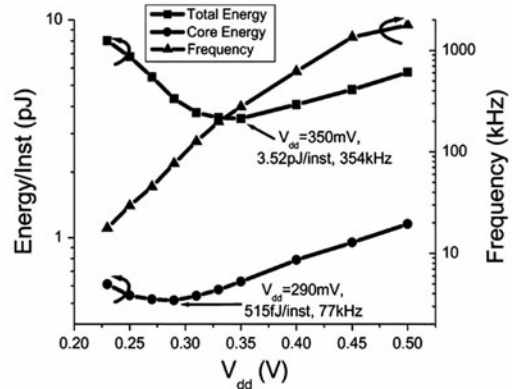
Figure 4.20a–d shows the energy trend and the presence of the MEP in various integrated prototypes, including an FFT core from MIT (Wang and Chandrakasan 2005), an 8-bit micro-processor from Umich (Hanson et al. 2008), an IA-32 processor from Intel (Jain et al. 2012), and an AES core from NUS (Zhao et al. 2015). From these figures, the energy curve is relatively flat around the MEP, hence the minimum- or nearly-minimum energy per operation does not require a stringent precision in the generation of the supply voltage. In practical designs, a change in V_{DD}

⁷ Since the energy per operation E_{TOT} in (5.11) is proportional to E_{cycle} , in the following we will simply refer to the energy per cycle in (5.12), unless otherwise specified. All considerations are immediately extended to E_{TOT} by simply multiplying E_{cycle} by CPO .

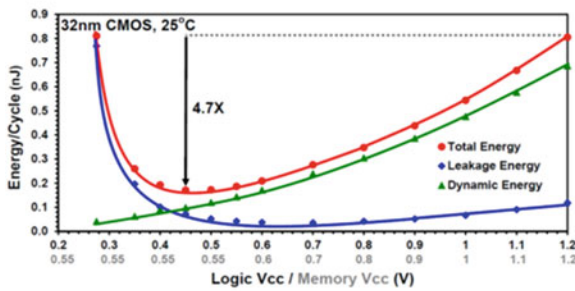
Fig. 4.19 Qualitative trend of dynamic, leakage and total energy per cycle E_{cycle} (or equivalently total energy per operation E_{TOT}) vs. supply voltage V_{DD}



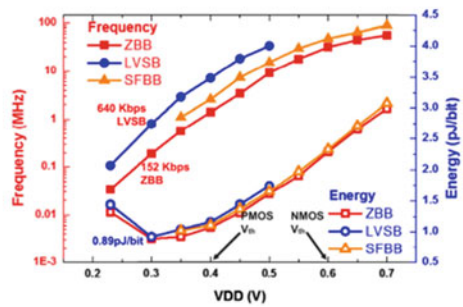
(a)



(b)



(c)



(d)

Fig. 4.20 Energy vs. V_{DD} and minimum-energy point in (a) FFT core ((Wang and Chandrakasan 2005) from MIT), (b) 8-bit microprocessor ((Hanson et al. 2008)

from Umich), (c) IA-32 processor ((Jain et al. 2012) from Intel), (d) AES core ((Zhao et al. 2015) from NUS)

around the optimal voltage $V_{DD,opt}$ by various tens of mVs (e.g., 30-50 mV) keeps the energy very close to E_{min} (e.g., within a few percentage

points). The MEP voltage $V_{DD,opt}$ in the above examples covers the typical range encountered in real designs (300–450 mV). The detailed

analysis on the dependence of the MEP position on process and design parameters is presented in the next subsection.

Let us observe that the leakage energy increase at low voltages limits the energy reductions enabled by aggressive voltage scaling, compared to the quadratic reduction that would be achievable if the total energy were dominated by E_{dyn} . Indeed, in the latter case the minimum achievable energy would be given by (4.9) with V_{DD} equal to the minimum operating voltage V_{min} that ensure correct operation, as in Fig. 4.19. The related energy saving compared to nominal voltage is reported in Table 4.2, which represents an upper bound of the energy savings achievable for quick estimates. Observe that the potential energy savings in circuits with wire-dominated load are lower than the case of gate-dominated load. Indeed, in the former case the dynamic energy reduction is purely quadratic, whereas the latter also benefits from the simultaneous load reduction due to the reduction in the transistor gate capacitance at low V_{DD} (see Sect. 4.2.1 and Fig. 4.4). From this table, operation at near-threshold voltages potentially reduces the energy by up to an order of magnitude, compared to the nominal voltage. At the same time, the presence of leakage narrows down the range of voltages at which energy reduction is truly allowed.

Table 4.2 Dynamic energy reduction vs. V_{DD}

V_{DD} (mV)	V_{DD}^2 energy saving (load = wire only)	$C_g \cdot V_{DD}^2$ energy saving (load = transistor only)
200 mV	36X	54X
400 mV	9X	11.6X
600 mV	4X	4.4X
800 mV	2.2X	2.4X
1 V	1.4X	1.4X
1.2 V	1X	1X

Table 4.3 Measured energy breakdown in Jain et al. (2012)

	V_{DD} (V)		$E_{lkg}(\%)$			$E_{dyn}(\%)$		
	Logic	L1C	Logic	L1C	Total	Logic	L1C	Total
Sub threshold (V_{min})	0.28	0.55	62%	33%	95%	4%	1%	5%
Near threshold (MEP, $V_{DD,opt}$)	0.45	0.55	27%	15%	42%	53%	5%	58%
Above threshold (nominal V_{DD})	1.2	1.2	11%	3%	14%	81%	5%	86%

logic = Core, L1C = L1 Cache

As even more crucial observation, the leakage energy is a substantially larger fraction of the overall energy budget at near-threshold voltages, compared to nominal voltage. Indeed, E_{lkg} (E_{dyn}) at near-threshold voltages is larger (smaller) than at nominal V_{DD} . Table 4.3 shows the detailed energy breakdown measured in the microprocessor in (Jain et al. 2012), which includes a level-1 cache. Above threshold, the leakage energy is 14% and is well in line with the expectations at nominal voltage. In near-threshold region, the leakage energy raises to a much larger 42% as expected, and in sub-threshold region it completely dominates the overall energy.

For all the above reasons, mitigating the leakage energy is a crucial goal of near-threshold designs, and is far more important than traditional low-power above-threshold designs. In addition, Sect. 4.4 will show that traditional low-power techniques to mitigate leakage are rather ineffective when V_{DD} is pushed down to near threshold.

4.3.3 Trans-Regional Energy Model

From the previous subsection, the MEP is set by the optimal balance between dynamic and leakage energy. In other words, the MEP voltage $V_{DD,opt}$ and the resulting minimum energy E_{min} both depend on the ratio between leakage and total energy. This means that the MEP position in the energy-voltage plane changes according to this ratio, as discussed below.

When the leakage energy significantly increases for some reason, whereas the dynamic energy remains constant, the total energy clearly increases and $V_{DD,opt}$ increases as well (i.e., the MEP moves to the right, and upwards, as summarized in Fig. 4.21a). Indeed, in this case

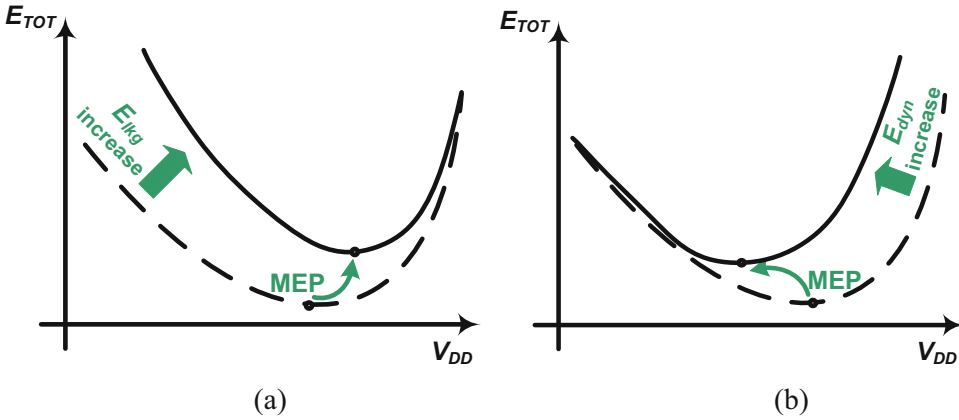


Fig. 4.21 Qualitative description of how the minimum-energy point (MEP) changes position when (a) E_{lkg} increases (E_{dyn} kept constant), (b) E_{dyn} increases (E_{lkg} kept constant)

the leakage energy tends to be a larger fraction of E_{TOT} , hence V_{DD} needs to be increased to reduce E_{lkg} (as explained by Fig. 4.19). On the other hand, when the dynamic energy increases at iso-leakage energy, E_{dyn} becomes a larger fraction of E_{TOT} , hence it becomes more important to reduce E_{dyn} and hence $V_{DD,opt}$ decreases (i.e., the MEP moves to the left, and upwards, as summarized in Fig. 4.21a). From the above considerations, the MEP tends to move to the right when the temperature is increased and/or the circuit activity is reduced, due to a different input data profile or power mode (i.e., different

modules are activated). Since the input dataset and the temperature are time varying and are unpredictable at design time, a feedback scheme that tracks the actual MEP through energy sensing (or estimation) and adjusts the supply voltage accordingly (Ramadass and Chandrakasan 2008).

To have a more quantitative understanding of the dependence of the MEP position on process, design and environmental parameters, let us consider an analytical model of the energy. In detail, eq. (4.12) can be written in the following more useful form⁸:

$$\begin{aligned}
 E_{cycle} &= \alpha_{SW} \cdot C_{TOT} \cdot V_{DD}^2 \left[1 + LD_{eff} \cdot \frac{I_{off,TOT}}{\alpha_{SW} \cdot C_{TOT} \cdot V_{DD}} FO4 \right] \\
 &\approx \alpha_{SW} \cdot C_{TOT} \cdot V_{DD}^2 \left[1 + LD_{eff} \cdot \frac{(\text{gatecount} \cdot \overline{I_{off}})}{\alpha_{SW} \cdot (\text{gatecount} \cdot C_{cell})} \cdot \frac{5C_{in,min}}{2I_{on,min}} \right] \\
 &\approx \alpha_{SW} \cdot C_{TOT} \cdot V_{DD}^2 \left[1 + 2.5 \cdot LD_{eff} \cdot \frac{I_{0,min} \cdot e^{-\frac{V_{TH}}{n \cdot \frac{kT}{q}}} \cdot \frac{\text{strength}}{X_{stack,off}}}{\alpha_{SW} \cdot \frac{C_{cell}}{C_{in,min}}} \cdot \frac{1}{I_{0,min} \ln \left(e^{\frac{V_{DD}-V_{TH}}{n \cdot (kT/q)}} + 1 \right)} \right] \\
 &= \alpha_{SW} \cdot C_{TOT} \cdot V_{DD}^2 \left[1 + ILDR \cdot e^{-\frac{V_{DD} \left(1 + \alpha_{X_{off}} \right)}{n \cdot \frac{kT}{q}}} \cdot f_{ILDR}(V_{DD}) \right] \quad (4.13)
 \end{aligned}$$

where (4.7) and (4.26) were used, and the intrinsic leakage-dynamic energy ratio $ILDR$ (i.e., the contribution of E_{lkg}/E_{dyn} that is independent of V_{DD}) was defined as

⁸ The following analysis is inspired by Hanson et al. (2006a), Hanson et al. (2006b), Bo et al. (2004) and generalizes the results to arbitrary designs, instead of being valid only for simple cascaded inverters.

$$ILDR = 2.5 \cdot LD_{eff} \cdot \frac{\overline{strength}}{X_{stack,off} \Big|_{V_{DD,nom}}} \cdot e^{-\alpha_{X_{off}} \cdot \frac{V_{DD,nom}}{n \cdot kT/q}} \cdot \frac{1}{\alpha_{SW} \cdot \frac{\overline{C_{cell}}}{C_{in,min}}} \quad (4.14a)$$

and $f_{ILDR}(V_{DD})$ was defined as

$$f_{ILDR}(V_{DD}) = \frac{e^{\frac{V_{DD}(1+\lambda_{DIBL})-V_{TH0}}{n \cdot kT/q}}}{\ln\left(e^{\frac{V_{DD}(1+\lambda_{DIBL})-V_{TH0}}{n \cdot kT/q}} + 1\right)} \quad (4.14b)$$

Also, in (4.13) it was observed that

- the total leakage current $I_{off,TOT}$ of the design under consideration is equal to the gate count (*gatecount*) multiplied by the average leakage per standard cell $\overline{I_{off}}$
- $\overline{I_{off}}$ can be expressed as the leakage current of a minimum-sized inverter $I_{0,min} \cdot e^{-V_{TH}/n \cdot kT/q}$, multiplied by the average cell strength $\overline{strength}$ (where 1X refers to the minimum-sized inverter) and divided by the average off-stacking factor $X_{stack,off}$ (i.e., the factor by which the leakage current is reduced due to transistor stacking)
- the total capacitance C_{TOT} is equal to the gate count multiplied by the average switched capacitance per standard cell $\overline{C_{cell}}$
- $FO4$ can be thought of as the delay of a minimum-sized inverter driving four equal inverters, and hence its total load capacitance is $4C_{in,min}$ ($C_{in,min}$ is the input capacitance of a minimum-sized inverter) plus its parasitic capacitance, which is approximately equal to $C_{in,min}$ (Sutherland et al. 1999); $I_{on,min}$ is the current delivered by such minimum-sized inverter (see (4.4)), and $I_{0,min}$ is the I_0 parameter in (4.1) and (4.2) pertaining to the same inverter.

In (4.13), the average off-stacking factor $X_{stack,off}$ is evaluated at the nominal voltage

$V_{DD,nom}$ of the adopted technology, and its down-scaling at low voltages is accounted for by the technology-dependent parameter $\alpha_{X_{off}} \ll 1$ (see (4.26)). Eq. (4.13) is strongly affected by $ILDR$. From (4.13), the latter parameter represents the voltage-independent (i.e., intrinsic) contribution of the ratio between leakage and dynamic energy, and is defined by

- the architecture through LD_{eff} (i.e., faster designs with low LD_{eff} exhibit lower $ILDR$)
- the function implemented by the design under consideration, which in turn sets the standard cell usage statistics (i.e., the average off-stacking factor $X_{stack,off}$) and the average fan-out $\overline{C_{cell}}/C_{in,min}$ (i.e., the average equivalent number of minimum-sized inverters that load the cells); such dependence tends to be fairly weak, due to the averaging effect across cells in large designs
- the performance target in the automated cell sizing phase, as $\overline{strength}$ and the average fan-out $\overline{C_{cell}}/C_{in,min}$ both tend to be larger for tighter timing constraints (again, faster designs have lower $ILDR$); accordingly, such dependence tends to be fairly weak as well
- the input dataset, which in turn sets the circuit activity (i.e., α_{SW}).

In summary, parameter $ILDR$ is mainly set by the architecture and the input statistics, and low values of $ILDR$ are associated with more active and faster designs. In practical cases, $ILDR$ ranges from a few hundreds for rather fast and active designs with heavy dynamic energy, to a several tens of thousands in very slow and inactive circuits. Higher values are observed only when an additional constant power contribution comes from external blocks.

Table 4.4 Numerical Examples for $ILDR$ in 28 nm ($V_{DD,nom} = 1.2$ V)

Design	LD_{eff}	α_{sw}	$\frac{\overline{C_{cell}}}{C_{in,min}}$	$\overline{strength}$	$ILDR$	$V_{DD,opt}$	$\frac{E_{lk}}{E_{dyn}} \Big _{MEP}$
High performance, active	20	15%	20	6	110	180 mV	0.62
Low performance, little active	100	3%	6	3	3100	320 mV	0.27

Observe that $ILDR$ is defined at nominal voltage and all parameters not explicitly related to $V_{DD,nom}$ are essentially independent of the voltage,⁹ and hence can be evaluated from the report of synthesis/place&route at such voltage without requiring the full characterization of the library at different voltages. A few numerical examples in 28 nm are reported in Table 4.4, assuming $\lambda_{DIBL} = 0.1$ (i.e., $\alpha_{X,off} = 0.098$ from (4.26)), $X_{stack,off} \Big|_{V_{DD,nom}}$ equal to 20 (i.e., average of two stacked transistors in this technology) at nominal voltage. From the technology scaling viewpoint, k_0 tends to slightly decrease at finer technologies, due to stronger DIBL and hence larger $X_{stack,off} \Big|_{V_{DD,nom}}$. As a simpler approach, $ILDR$ can also be estimated as the value that makes the ratio E_{lkq}/E_{dyn} (i.e., $ILDR \cdot e^{-V_{DD}(1+\alpha_{X,off})/n \frac{kT}{q}}$ $f_{ILDR}(V_{DD})$ in (4.13)) equal to the value that is obtained from power analysis at RTL level.

4.3.4 Considerations on the MEP Voltage

Typically, the MEP mostly lies in the deep sub-threshold region (Hanson et al. 2006a), and sometimes near-threshold (Hanson et al. 2006b). In the former case, $f(V_{DD}) \approx 1$ in (4.13) since $V_{DD} < V_{TH0}$ in (4.14b), hence the energy is independent of the transistor threshold voltage. This is because the latter affects both the leakage and the on-current in the same way (i.e., both are proportional to $\exp\left(-V_{TH}/n \cdot \frac{kT}{q}\right)$ in sub-threshold). In this case, V_{TH} is chosen

⁹For example, the average fan-out $\overline{C_{cell}}/C_{in,min}$ is independent of V_{DD} since the wire capacitance is constant, and the transistor gate capacitance does not change substantially (see Fig. 5.4). Similarly, the logic depth, the activity and the average strength do not depend on V_{DD} .

exclusively based on the performance requirement (i.e., targeted $FO4$), according to (4.4).

Let us now analyze the optimum voltage $V_{DD,opt}$ that minimizes the energy in (4.13) assuming the MEP to be in sub-threshold. Although a closed-form solution cannot be found, a good approximation for a single- V_{TH} design is logarithmic (similar to (Bo et al. 2004; Hanson et al. 2006a, b))

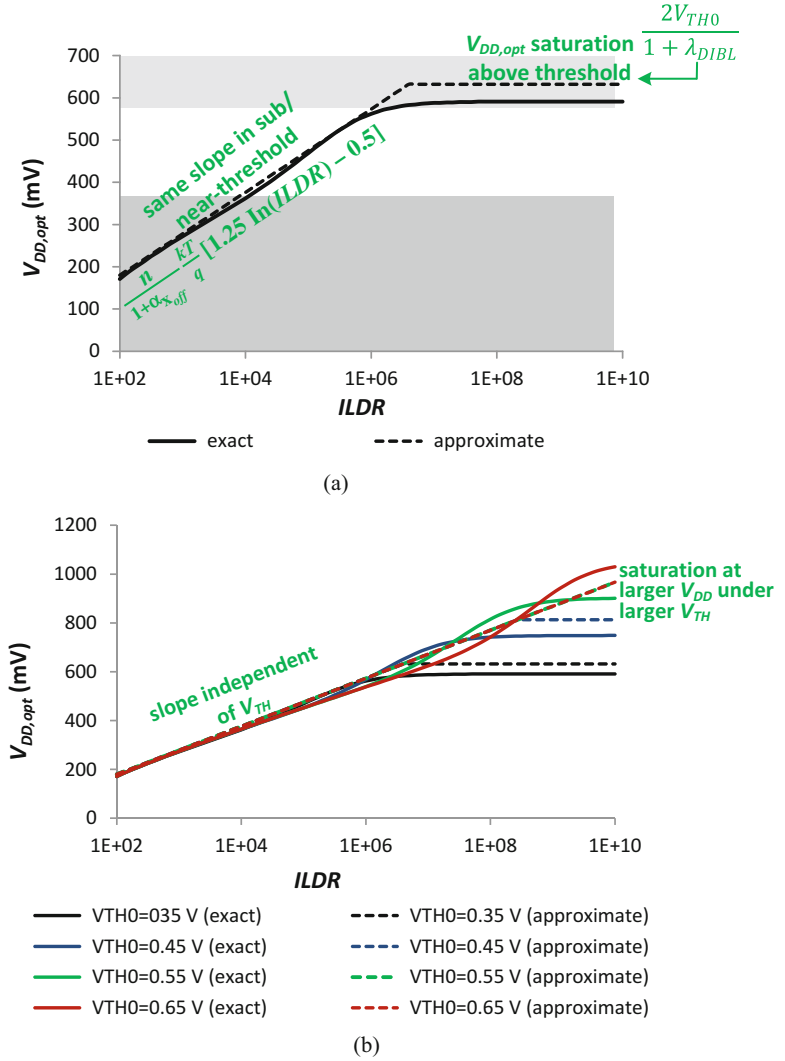
$$V_{DD,opt} \approx \frac{n}{1 + \alpha_{X,off}} \frac{kT}{q} [1.25 \ln(ILDR) - 0.5] \quad (4.15)$$

which has a maximum error of 4% for practical values of $ILDR$, as plotted in Fig. 4.22a under the above 28-nm parameters and $V_{TH0} = 0.35$ V. From this figure, the MEP voltage logarithmically increases with the constant slope in (4.15), which is independent of the adopted threshold voltage, as shown in Fig. 4.22b.

As expected from the above considerations and Fig. 4.21, the MEP moves to the right for slow and less active circuits, and to the left for circuits with dominating dynamic energy and low logic depth. Observe that operation at $V_{DD} < V_{min}$ severely degrades the die yield, hence minimum-energy designs need to adopt a supply voltage equal to the minimum between $V_{DD,opt}$ in (4.15) and V_{min} . From Fig. 4.22a, b, this means that fast (i.e., with low LD_{eff}) designs with $ILDR$ lower than a few thousands cannot really achieve true-minimum energy, due to voltage scaling limitations imposed by robustness issues.

The above analysis assumed that the MEP lies in the deep sub-threshold region, which is correct as long as $V_{DD,opt}$ in (4.13) is lower than $V_{TH} - 50$ mV (see Fig. 4.1), i.e., when $ILDR < e^{\frac{V_{TH}-50 \text{ mV}}{1.5n \frac{kT}{q}} + 1.6}$. The latter boundary value for $ILDR$ in 28 nm is typically in the order of 1,000–2,000 for $V_{TH} = 350$ mV (including

Fig. 4.22 MEP voltage vs $ILDR$ for 28-nm technology with (a) $V_{TH0} = 0.35\text{ V}$ and detailed analytical model, (b) $V_{TH0} = 0.35\text{ V}, 0.45\text{ V}, 0.55\text{ V}, 0.65\text{ V}$ (exact solution via numerical minimization of (4.13), approximate expression as in (4.15)–(4.17))



DIBL), 4000–5000 for $V_{TH} = 400\text{ mV}$, and 10,000 for $V_{TH} = 450\text{ mV}$. Slightly larger values are typically found in older technologies, due to the lower subthreshold factor n .

Interestingly, Fig. 4.22a, b show that (4.15) can be extended to the near-threshold region (i.e., larger $ILDR$), as it still predicts the MEP voltage with good accuracy. Hence, the MEP voltage is again independent of the threshold voltage, even in the near-threshold region. For even larger values of $ILDR$ such that $V_{DD} > V_{TH0} + 200\text{ mV}$ (see Fig. 4.1), the MEP moves to the above-threshold region and eventually saturates to a

value $V_{DD,opt}|_{ILDR \rightarrow \infty}$. Indeed, $f(V_{DD})$ in (4.14b) becomes approximately equal to $e^{\frac{V_{DD}(1+\lambda_{DIBL})-V_{TH0}}{n \frac{kT}{q}}} \cdot \left(\frac{V_{DD}(1+\lambda_{DIBL})-V_{TH0}}{n \cdot (kT/q)}\right)$, and the resulting $V_{DD,opt}$ is found by minimizing (4.13) for $ILDR \rightarrow \infty$:

$$V_{DD,opt}|_{ILDR \rightarrow \infty} = \frac{2V_{TH0}}{1 + \lambda_{DIBL}} \quad (4.16)$$

which was found to be always within 12% of the exact solution that minimizes (4.13) in 28 nm (and typically within 5%). $V_{DD,opt}$ saturates because E_{lkg} increases when increasing V_{DD} in

the above-threshold region, as opposed to sub- and near-threshold (see Fig. 4.18 and related discussion). In other words, it does not make sense to increase V_{DD} beyond (4.16) from an energy viewpoint, as this would surely increase both dynamic and leakage energy, and hence the total energy. Indeed, (4.16) represents the

$$V_{DD,opt} \approx \min\left(\frac{n}{1 + \alpha_{X_{off}}}\frac{kT}{q}[1.25\ln(ILDR) - 0.5], \frac{2V_{TH0}}{1 + \lambda_{DIBL}}\right) \quad (4.17)$$

which has a typical (maximum) error of 4% (7%) across the very wide range of $ILDR$ in Fig. 4.22a, b. Eq. (4.17) is a useful tool to estimate the MEP position by knowing the type of design (i.e., $ILDR$), and a few other technology-dependent parameters. From (4.17), the transistor threshold choice affects only the value of $ILDR$ and the voltage at which the MEP saturates at. In particular, larger V_{TH0} moves saturation towards exponentially larger k_0 and proportionally larger $V_{DD,opt}|_{ILDR \rightarrow \infty}$.

voltage at which E_{lkg} is minimum (i.e., such that $V_{DD} \cdot I_{off} \cdot FO4$ is minimum, from (4.12)).

In summary, the above considerations suggest that $V_{DD,opt}$ can be simply modeled by extending (4.15) to the near-threshold and part of the above-threshold region, and limiting it to its asymptotic maximum value in (4.16):

4.3.5 Considerations on the MEP Energy

From (4.13), the resulting energy at the MEP in deep sub-threshold region (i.e., under (4.15)) can be written as

$$E_{min} = \alpha_{SW} \cdot C_{TOT} \cdot V_{DD,opt}^2 \left[1 + \frac{E_{lkg}}{E_{dyn}}\Big|_{MEP}\right] \quad (4.18)$$

where, considering that $f(V_{DD}) \approx 1$ and $\alpha_{X_{off}} \ll 1$ in (4.13), the energy-optimum ratio E_{lkg}/E_{dyn} at the MEP is given by

$$\frac{E_{lk}}{E_{dyn}}\Big|_{MEP, sub-threshold} \approx ILDR \cdot e^{-\frac{V_{DD,opt}}{n\frac{kT}{q}}} \approx 0.2 + \frac{17}{ILDR^{0.75}} \quad (4.19)$$

In (4.19), an empirical approximate expression has been introduced to facilitate its estimate at design time, and its error is within 12% for low values of $ILDR$ (down to 150), as plotted in Fig. 4.23. Observe that E_{lkg}/E_{dyn} at the MEP is independent from the chosen (single) threshold voltage, as expected from the considerations in Sect. 4.3.4. In other words, sub-threshold designs that differ only for the threshold voltage choice have the same leakage percentage contribution, other than the same $V_{DD,opt}$ (see Sect. 4.3.4). Accordingly, the MEP is hence chosen based on the performance target rather than energy. Also, this means that the MEP voltage can be estimated at design time even before choosing the transistor flavor (and hence before actual implementation).

Compared to the overall energy budget, E_{lkg} at the MEP needs to be 40–50% for very fast/active designs ($ILDR \leq 150$), around 15–30% for more typical designs ($ILDR > 150$ but still in sub-threshold). Previous work on joint supply/threshold voltage and sizing optimization showed that energy optimality is achieved when E_{lkg} is about one third of the overall energy (Markovic et al. 2004; Nose et al. 2000; Patil). Accordingly, these results hold in sub-threshold region only for relatively slow and inactive designs, from Fig. 4.23.

As discussed in Sect. 4.3.4, very fast and active designs have low $V_{DD,opt}$, which often times falls below V_{min} . In these cases, minimum-energy and reliable operation is

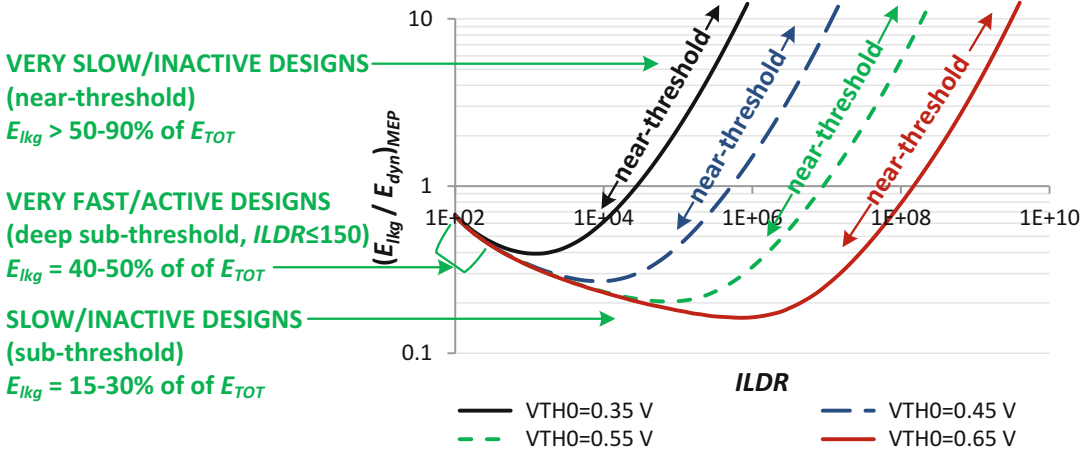


Fig. 4.23 Leakage-dynamic energy ratio at the MEP vs. $ILDR$, obtained through numerical minimization of (4.13a)

achieved at $V_{DD} = V_{min} > V_{DD,opt}$. If the extra performance compared to the MEP is not utilized since the design is essentially energy constrained, operation at $V_{min} > V_{DD,opt}$ leads to an increase in E_{dyn} by a factor $(V_{min}/V_{DD,opt})^2$ from (4.9), compared to the MEP. At the same time, a smaller increase in E_{lkg} by a factor $V_{min}/V_{DD,opt}$ is observed (since same clock cycle is assumed in (4.10), and DIBL effect is neglected). In other words, designs with $V_{min} > V_{DD,opt}$ typically have lower E_{lkg}/E_{dyn} , compared to operation at the MEP in Fig. 4.23.

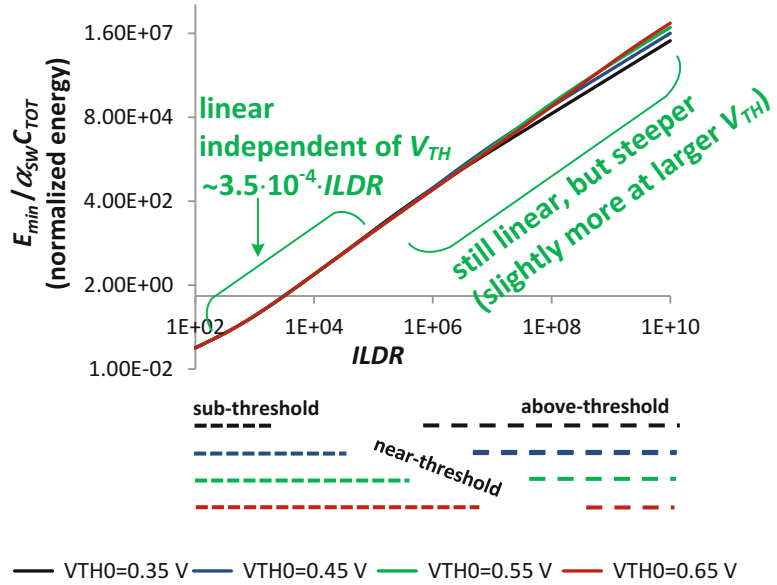
For larger $ILDR$, again the MEP moves to the near-threshold region and the energy-optimum ratio E_{lkg}/E_{dyn} at the MEP increases again when increasing V_{DD} . This is because the increase in $V_{DD,opt}$ at near-threshold voltages determines a much smaller reduction in E_{lkg} compared to sub-threshold, as the gate delay decreases much slower than exponentially from (4.6). Analytically, this is accounted for by the increase in $f(V_{DD})$ in (4.13), which determines a proportional increase in E_{lkg}/E_{dyn} . Accordingly, E_{lkg} becomes again a substantial fraction of the energy budget when the MEP is pushed at near-threshold voltages or higher (i.e., for large $ILDR$), as shown in Fig. 4.23. This explains why E_{lkg}/E_{dyn} heavily depends on V_{TH} at near-threshold voltages, as in Fig. 4.23. For extremely slow and inactive circuits, the MEP moves to the

above-threshold region, and is definitely dominated by E_{lkg} .

From the above considerations, the energy E_{min} at the MEP monotonically increases when increasing $ILDR$. At sub-threshold voltages, this is due to the increase in the energy-optimal voltage $V_{DD,opt}$ in (4.17), which is certainly more rapid than the reduction in $(1 + E_{lkg}/E_{dyn}|_{MEP})$ in (4.19). Since both dependencies were found to be unaffected by the threshold voltage in sub-threshold, E_{min} for MEP in sub-threshold is independent of V_{TH} as well. At near-threshold voltages, E_{min} keeps increasing since $V_{DD,opt}$ in (4.17) continues to increase with the same trend as sub-threshold (see Fig. 4.22a), and $(1 + E_{lkg}/E_{dyn}|_{MEP})$ increases as well. For above-threshold MEP, $V_{DD,opt}$ in (4.17) saturates to an almost constant value, and $(1 + E_{lkg}/E_{dyn}|_{MEP})$ keeps increasing.

From the above considerations, the energy at the MEP is monotonically degraded when $ILDR$ is increased, i.e., for leaky or little active designs. This is essentially due to the increase in $V_{DD,opt}$ (i.e., dynamic energy at the MEP), and the increase in the leakage-dynamic energy ratio at voltages above V_{TH} . More quantitatively, Fig. 4.24 shows that E_{min} increases in an approximately linear fashion when increasing $ILDR$. In

Fig. 4.24 Energy at MEP E_{min} normalized to $\alpha_{sw}C_{TOT}$ vs. $ILDR$ for various threshold voltages



particular, in the sub-threshold region the ratio $E_{min}/\alpha_{sw}C_{TOT}$ is well approximated by $3.5 \cdot 10^{-4} \cdot ILDR$ regardless of the threshold voltage, hence

$$E_{min} \approx 3.5 \cdot 10^{-4} \cdot \alpha_{sw} \cdot C_{TOT} \cdot ILDR \quad (4.20)$$

which is within 20% of exact E_{min} for $ILDR$ up to a few tens of thousands. For less typical design with larger $ILDR$, the trend becomes slightly steeper by a factor ranging from to 2.5X to 4X compared to (4.20), for a threshold voltage in the 350–650 mV range. In other words, when the MEP is at near-threshold voltages, E_{min} actually increases when V_{TH} increases, although moderately.

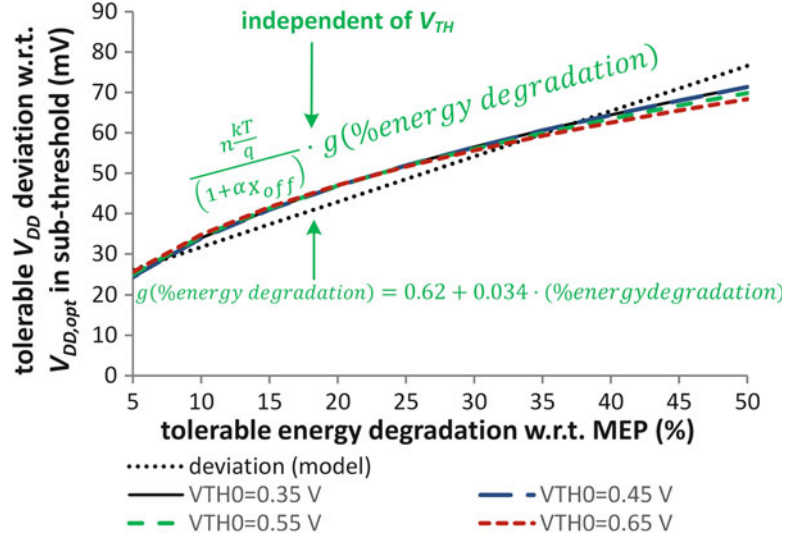
Summarizing these conclusions and those in Sect. 4.3.4, the MEP voltage is unaffected by the (single) threshold voltage when operating in sub- and near-threshold voltages. On the other hand, the energy portion associated with leakage drastically increases when operating at near-threshold voltages, compared to sub-threshold ones. At above-threshold voltages, the MEP voltage becomes a function of V_{TH} , and energy is dominated by leakage. Regardless of the voltage range in which the MEP lies in, the minimum achievable energy monotonically and proportionally increases with $ILDR$.

4.3.6 Sensitivity of Nearly-Minimum Energy to V_{DD} Inaccuracies

From a design standpoint, it is necessary to predict the required accuracy for V_{DD} to achieve nearly-minimum energy per operation, which in turn constraints the design of the voltage regulation circuitry and the power management sub-system. As can be seen from Fig. 4.20a–d, the energy-voltage curve is steeper at the left of the MEP, due to the exponential increase in the leakage energy at low voltages. In other words, an uncertainty $\pm \Delta V_{DD}$ in the supply voltage around the MEP degrades the energy more substantially when it pushes V_{DD} below $V_{DD,opt}$ rather than above (even more so, if performance is considered). Due to the same reason, the energy degradation at the left of the MEP compared to the right becomes more evident for larger ΔV_{DD} .

In nearly-minimum energy designs, the maximum tolerable percentage energy degradation $\% \text{energydegradation}$ compared to the MEP due to the uncertainty in V_{DD} needs to be translated into the specification of the maximum tolerable uncertainty ΔV_{DD} . In sub-threshold region, the resulting tolerable uncertainty ΔV_{DD}

Fig. 4.25 Maximum supply voltage deviation from MEP that maintains the energy degradation within the target (on x-axis) in sub-threshold region (model in (4.21a), (4.21b))



is independent of V_{TH} , since the energy is independent of V_{TH} as well (see Sect. 4.3.4). Since $f(V_{DD}) \approx 1$, Eq. (4.13) can be expressed in a technology-independent manner¹⁰ by defining the normalized voltage $V_{DD,norm} = V_{DD}(1 + \alpha_{X_{off}})/n \frac{kT}{q}$. The maximum deviation $\Delta V_{DD,norm}$ in $V_{DD,norm}$ compared to the value that minimizes the energy can be easily solved numerically. The numerical

solution $\Delta V_{DD,norm}$ turns out to be largely independent of $ILDR$, and is hence only a function of % energy degradation. $\Delta V_{DD,norm}$ is well approximated (within 10%) by $0.62 + 0.034 \cdot (\%energydegradation)$, as shown in Fig. 4.25. Accordingly, the maximum tolerable voltage deviation that meets a targeted percentage deviation in sub-threshold region is

$$\Delta V_{DD,subthreshold} \approx \frac{n \frac{kT}{q}}{(1 + \alpha_{X_{off}})} \cdot g(\%energydegradation) \quad (4.21a)$$

$$g(x) = 0.62 + 0.034 \cdot x \quad (4.21b)$$

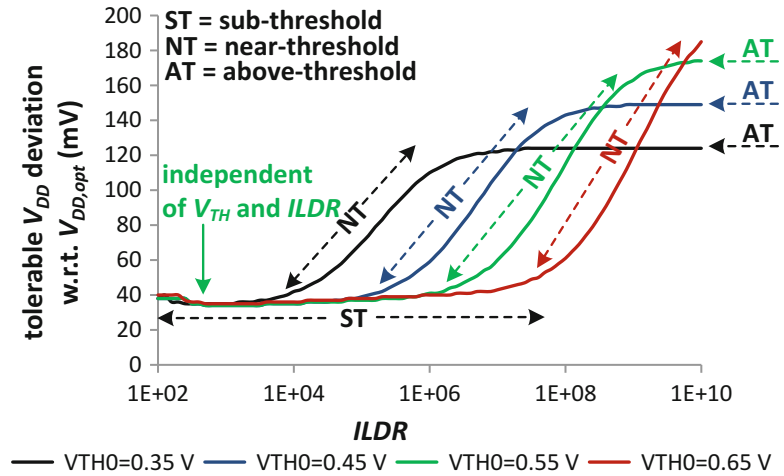
Interestingly, from (4.21a), (4.21b), ΔV_{DD} in sub-threshold does not depend on the position of the MEP, and it only depends on technology through subthreshold slope and DIBL coefficient, and on the targeted maximum energy degradation. As an example, Fig. 4.25 plots the maximum tolerable ΔV_{DD} versus the energy degradation in 28 nm, and shows that V_{DD} needs

to be set with a precision of about a thermal voltage (25–35 mV) to keep the energy degradation compared to the MEP modest (5–10%). Larger V_{DD} uncertainty (e.g., 1.5–2X the thermal voltage) leads instead to an unacceptably large energy degradation, and should hence be avoided in practical cases.

When the MEP moves to the near-threshold region, a larger voltage deviation can be tolerated for a targeted maximum energy degradation compared to the MEP. This is because $FO4$ and hence E_{lkg} (see Eq. (4.10)) become less sensitive to V_{DD} compared to sub-threshold, as shown in Fig. 4.26. As expected, the tolerable voltage

¹⁰ Indeed, eq. (5.13) in sub-threshold region becomes $E_{cycle} \propto V_{DD,norm}^2 [1 + ILDR \cdot e^{-V_{DD,norm}}]$, assuming $\alpha_{X_{off}} \approx \lambda_{DIBL}$ (which is generally, since $\lambda_{DIBL} \ll 1$).

Fig. 4.26 Maximum supply voltage deviation from MEP vs. $ILDR$ for different V_{TH0} in 28 nm (target energy degradation w.r.t. MEP = 10%)



deviation at near-threshold depends on V_{TH0} , as opposed to sub-threshold. This is because the energy in (4.13) depends on V_{TH0} at near threshold (see Sect. 4.3.4), since V_{TH0} defines the voltage range (and hence $ILDR$) in which transistors enter this region. From Fig. 4.26, 4 to 6 thermal voltages can be tolerated with minimal energy penalty at near-threshold voltages.

For larger MEP voltages in the above-threshold region, an even larger voltage deviation around the MEP is tolerable for a given allowed energy degradation. This is because $FO4$ has the minimum sensitivity to V_{DD} across voltages from (4.6). As a consequence of the saturation of the MEP voltage discussed in Sect. 4.3.4, the maximum voltage

deviation saturates as well at above-threshold voltages, as shown in Fig. 4.26. As expected from (4.16), larger V_{TH0} pushes the saturation to higher voltages (and hence $ILDR$). Analytically, the maximum tolerable allowed voltage deviation around the MEP is evaluated by equating (4.13) and the energy at MEP (i.e., voltage in (4.16)) increased by a factor $(1 + \%energy\ degradation/100)$, and solving for V_{DD} . By approximating $\ln\left(e^{\frac{V_{DD}(1+\lambda_{DIBL})-V_{TH0}}{n(kT/q)}} + 1\right) \approx \frac{V_{DD}(1+\lambda_{DIBL})-V_{TH0}}{n(kT/q)}$ and $\alpha_{x_{off}} \approx \lambda_{DIBL}$ in (4.13), the maximum voltage deviation at above-threshold voltages results to

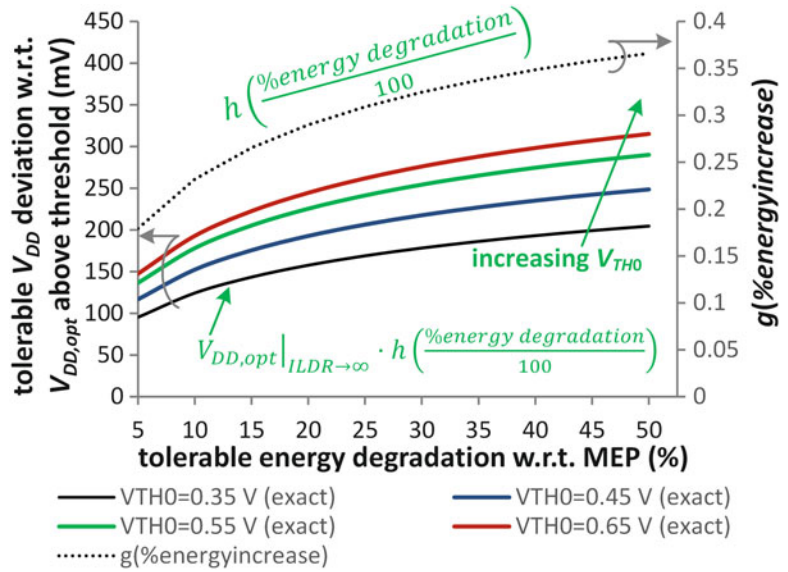
$$\Delta V_{DD,above-threshold} = V_{DD,opt}|_{ILDR \rightarrow \infty} \cdot h\left(\frac{\%energy\ degradation}{100}\right) \quad (4.22a)$$

$$h(x) = -x + \sqrt{x \cdot (1+x)} \quad (4.22b)$$

the first of which is plotted in Fig. 4.27 for a 28-nm technology, along with the technology-independent curve in (4.22b). Equation (4.22a), (4.22b) was found to be within 10-20% of the exact solution, for typical threshold voltages. For large V_{TH0} (e.g., 0.65 V), the error increases to 30-40% since ΔV_{DD} in (4.22a) becomes so large that it intrudes the near-threshold region, and the

above calculations hence become inaccurate. From (4.22a), (4.22b) the maximum voltage deviation around MEP above threshold depends on the technology through a proportional dependence on $V_{TH0}/(1 + \lambda_{DIBL})$. In other words, the maximum voltage deviation around the MEP above threshold is a fixed and technology-independent fraction of the MEP voltage, as set by $h(\% energy\ degradation/100)$. As opposed to sub-threshold region, the maximum deviation

Fig. 4.27 Maximum supply voltage deviation from MEP that maintains the energy degradation within the target (on x-axis) in above-threshold region (model in (4.22a), (4.22b))



around a MEP lying above threshold depends on V_{TH0} , and larger thresholds further relax the precision requirement on the voltage optimization and delivery. This is because a larger V_{TH0} enlarges the voltage range in which the MEP effectively increases for larger $ILDR$ (i.e., from about V_{TH0} up to (4.16)).

In summary, nearly-minimum energy operation requires the voltage to be controlled within approximately one thermal voltage when the MEP is in the sub-threshold region, independently of the position of the MEP. This requirement is substantially relaxed at near threshold, and increases at above threshold until saturation to a value that is proportional to the threshold voltage and sub-linearly related to the tolerable energy degradation. More in detail, deviation increases up to 4 thermal voltages for relatively low V_{TH0} , whereas it increases to more than 6 thermal voltages under large V_{TH0} .

Overall, this suggests that the performance can be increased with modest energy penalty by raising the voltage compared to the MEP, when the latter is at near- or above-threshold voltages. These considerations are summarized in the example in 28 nm in Fig. 4.28, which plots ΔV_{DD} versus $ILDR$ for various energy degradation targets and the related analytical models.

4.3.7 Example: ARM Core Operating at Minimum Energy

As a further numerical example, let us apply the above energy and MEP models to the design of the ARM Cortex M0 core in Myers et al. (2015). Table 4.5 reports all technology-, design- and workload-dependent parameters for this design, as obtained from the process design kit and information provided in Myers et al. (2015). The two programs “checksum” and “AES” are considered to consider a wide range of activities, from low (checksum) to high (AES), and hence observe the MEP shift due to different activity factors (the latter has 60% higher activity than the former (Myers et al. 2015)).

The resulting energy curve versus V_{DD} from experimental results in Myers et al. (2015) and the model in (4.13) is plotted in Fig. 4.29. Very good agreement can be observed across the wide voltage range from 0.25 to 1.2 V, with an average error of 1.7%, and an error well below 10% down to 0.3 V. As expected from Sect. 4.3.3 and Fig. 4.21, the MEP for the AES program is pushed to the left of the MEP for the checksum program, due to the higher activity determined by the former one. More quantitatively, the $ILDR$ factor in (4.13) and (4.14a), (4.14b) from the

Fig. 4.28 Summary of maximum V_{DD} deviation from MEP that maintains the energy degradation within the target vs. $ILDR$ and related models

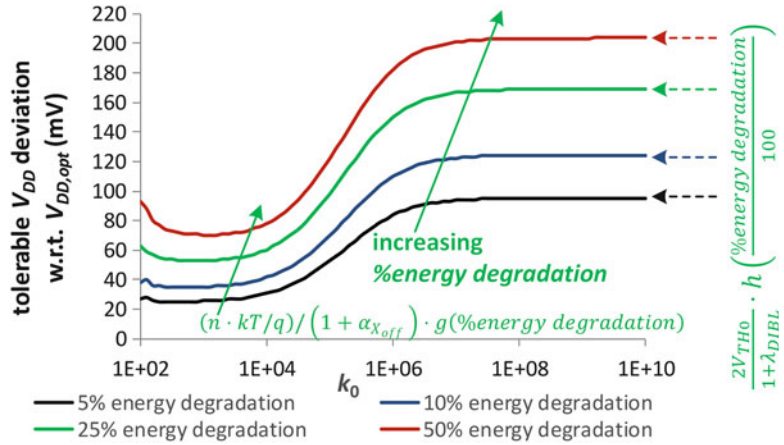


Table 4.5 Technology-, design- and workload-dependent parameters for ARM Cortex M0 Core in Myers et al. (2015)

	Parameter	Value
Technology	<i>Process</i>	65 nm
	$V_{DD,nom}$	1.2 V
	$V_{TH} @ V_{DD,nom}$	0.47 V
	λ_{DIBL}	0.095 V/V
	$\alpha_{x_{off}}$	0.087 V/V
	$FO4 @ V_{DD,nom}$	60 ps
Architecture/ckt design	LD_{eff}	240 ^a
	$\frac{C_{cell}}{C_{in,min}}$	4 ^b
	<i>strength</i>	2 ^b
Workload	α_{SW} (checksum program)	5% ^c
	α_{SW} (AES program)	8% ^c

^aEstimated from cycle time at nominal voltage (14.3 ns) and $FO4$ at nominal voltage

^bTypical values for very slow and low-energy designs (changes in a reasonable range do not significantly influence results)

^cActivity factor in AES obtained via a 60% increase compared to checksum program (Myers et al. 2015)

parameters in Table 4.5 results to 18,000 for the checksum program, and 11,200 for the AES program. The resulting voltage and energy at the MEP are summarized in Table 4.6, for both the experimental results in Myers et al. (2015) and the above models.

From Table 4.6, the estimated MEP voltage of the core in Myers et al. (2015) from Eq. (4.17) is 378 mV for the checksum program, which is close to the measured MEP voltage of 390 mV. The resulting minimum energy estimate of 11.6 pJ from (4.18) is also close to the measured energy of 11.7 pJ (Myers et al. 2015). At the MEP, the leakage energy is estimated to be smaller than the dynamic energy by a factor of 0.21 from eq. (4.19), which is close to the value of 0.22 in Myers et al. (2015). Good agreement of the models is also confirmed for the AES program, from the same table. Finally, the maximum tolerable V_{DD} uncertainty for a 10% energy increase compared to the MEP results to 45 mV from (4.21a), (4.21b), which agrees well with the value of approximately 48 mV in Myers et al. (2015).

4.4 Exploration of MEP Dependence on Logic Depth, V_{TH} Activity and Ineffectiveness of Leakage Reduction Techniques

In this section, the impact of logic depth, threshold voltage and activity are quantitatively and widely explored by considering the reference circuit in Fig. 4.30, applying the insights gained in Sect. 4.3. The simplicity and regularity of the circuit in 4.30 permits to gain an intuitive

Fig. 4.29 Experimental energy curve vs. V_{DD} in ARM Cortex M0 core (Myers et al. 2015) and energy predicted by the model in (4.13)

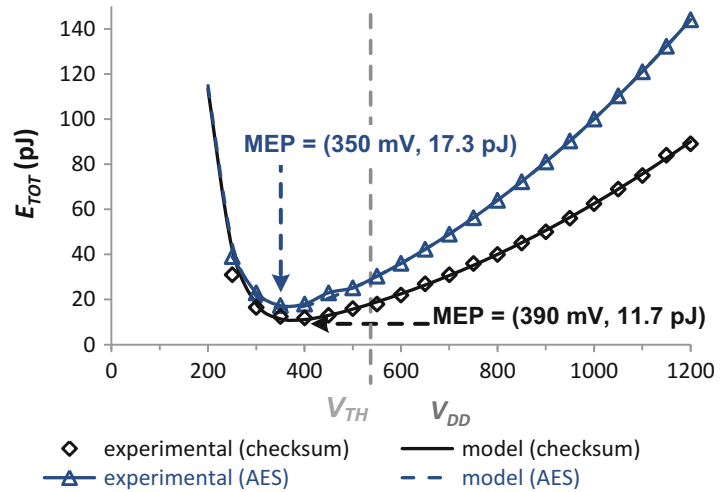


Table 4.6 Minimum Energy Point, Leakage/Dynamic Energy Ratio and Maximum Tolerable V_{DD} Uncertainty from (Myers et al. 2015) and Above Models

	checksum program		AES program	
	experimental (Myers et al. 2015)	model (equation)	experimental (Myers et al. 2015)	model (equation)
$V_{DD,opt}$	390 mV	378 mV (4.17)	350 mV	360 mV (4.17)
E_{min}	11.7 pJ	11.6 pJ (4.18)	17.3 pJ	16.8 pJ (4.18)
$\frac{E_{leak}}{E_{dyn}} _{MEP}$	0.22	0.21 (4.19)	0.22	0.22 (4.19)
ΔV_{DD} (energy increase = 10%)	45 mV	48 mV (4.21a), (4.21b)	45 mV	48 mV (4.21a), (4.21b)

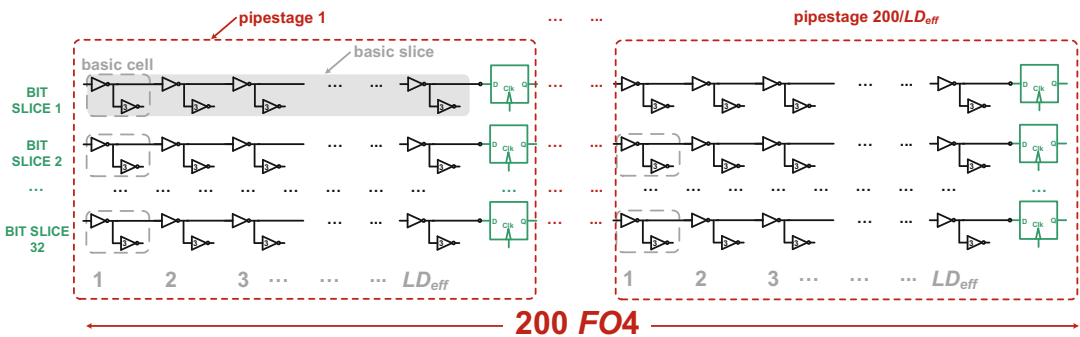


Fig. 4.30 Reference circuit for evaluation of the impact of logic depth, threshold voltage and activity

understanding of the underlying tradeoffs. Such circuit contains 32 slices of inverter gates, each with a fan-out of 4, and with a total logic depth LD_{TOT} (and hence delay by construction) of

$200FO_4$, as representative of a relatively complex microprocessor. The slices are interrupted through the insertion of registers, whose number is adjusted to achieve a targeted logic depth

LD_{eff} . Registers are made up of transmission-gate flip-flops, which are customarily encountered in standard cell libraries (Alioto et al. 2015).

4.4.1 Impact of Logic Depth

The heavy impact of E_{lkg} at near- and sub-threshold voltages can be mitigated by adopting microarchitectures with lower logic depth (i.e., deeper pipelining), from (4.10). However, deeper pipelining should be applied

judiciously to avoid a significant increase in the clocking overhead, which might offset some of the benefit brought by reduction in E_{lkg} . In the following, we will assume that the additional clocking cost of meeting the timing constraints with lower logic depth is modest (which is typically true in microarchitectures with $LD_{eff} \geq 25$ $FO4/cycle$).

The reference circuit in Fig. 4.30 has an overall energy per cycle equal to

$$E_{cycle} = \alpha_{SW} \cdot C_{TOT} \cdot V_{DD}^2 + \frac{LD_{TOT}}{LD_{eff}} \cdot E_{REG} + V_{DD} \cdot I_{off} \cdot FO4 \cdot LD_{eff} \quad (4.23)$$

where it was assumed that pipestages are perfectly balanced (i.e., the number of pipestages is $\frac{LD_{TOT}}{LD_{eff}}$). In (4.23), E_{REG} is the energy consumed by a single register, and $\frac{LD_{TOT}}{LD_{eff}}$ represents the

number of registers in the above circuit. From (4.23), an energy-optimal logic depth exists at a given V_{DD} , and its expression is readily found to be

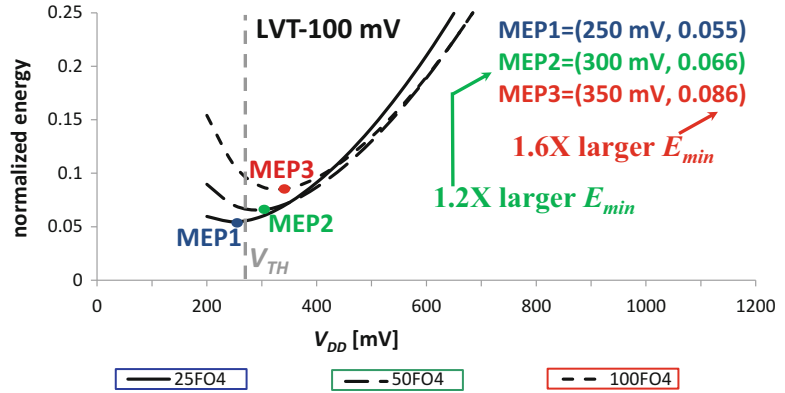
$$LD_{eff} = LD_{TOT} \sqrt{\frac{E_{REG}}{V_{DD} \cdot I_{off} \cdot LD_{TOT} \cdot FO4}} = LD_{TOT} \sqrt{\frac{E_{REG}}{E_{lkg}}} \quad (4.24)$$

where (4.10) was used to express E_{lkg} . From (4.24), the optimal logic depth is determined by the balance between the clocking and the leakage energy, as a larger number of registers leads to an increase in the former and a decrease in the latter. Such tradeoff is not really observed in traditional low-power (above-threshold) designs, as the leakage energy is usually kept a small fraction of the overall budget through several techniques (Narendra and Chandrakasan 2006), and the amount of pipelining is mainly defined by the performance target, or the dynamic energy-performance tradeoff under dynamic voltage scaling. Instead, nearly-minimum energy designs require a careful management of the clocking-leakage energy tradeoff, due to their strong interdependence (Alioto 2012). In addition, this means that energy-centric (micro)architectures need to be tailored around the targeted operation

voltage, and traditional architectures conceived for nominal voltage operation tend to be energy inefficient at low voltage. In other words, ultra-low power architectures need to be deeply rethought to truly enable nearly-minimum energy operation, as discussed in Chap. 3.

Quantitatively, eq. (4.24) suggests that the energy-optimal pipeline depth LD_{TOT}/LD_{eff} is given by the square root of the leakage-clocking energy ratio. Considering the large contribution of E_{lkg} at near-threshold voltages, the theoretical energy-optimal pipedepth tends to be quite small. In (4.24), the energy cost of all registers E_{REG} is assumed to be the same, since it refers to the simple reference circuit in Fig. 4.30. In more general architectures, the number of flip-flops per register, and hence the energy cost of a register, increases super-linearly under higher pipedepths (Chinnery and Keutzer 2007). Indeed,

Fig. 4.31 Energy normalized to value at nominal voltage vs. V_{DD} for logic depth of $25FO4$, $50FO4$ and $100FO4$



the overall number of flip-flops in a digital module increases according to a power law $(LD_{TOT}/LD_{eff})^{LGF}$, where $LGF > 1$ is the Latch Growth Factor, which is mainly defined by the specific function implemented (Srinivasan et al. 2002). Hence, the energy-optimal logic depth in general architectures tends to be moderately larger than predicted by (4.24).

On the low side, the energy-optimal logic depth is practically limited by the rapidly increasing clocking energy cost at small logic depths (i.e., deep pipelines). Typically, low logic depths in the order of 20–25 $FO4$ or smaller have a disproportionately large energy cost in the clock network at ultra-low voltages, and require non-straightforward clock network design approaches. The necessity of “fast” circuit and architectural designs¹¹ with deep pipelining at ultra-low voltages was first shown in Jeon et al. (2013), where an aggressive 17 $FO4$ logic depth was adopted in a 1024-point complex FFT processor. The adoption of such deep pipeline led to 17.7 nJ/transform at $V_{DD,opt} = 270$ mV in 65-nm CMOS, which was a 3.6X lower energy than previous state of the art. However, this required

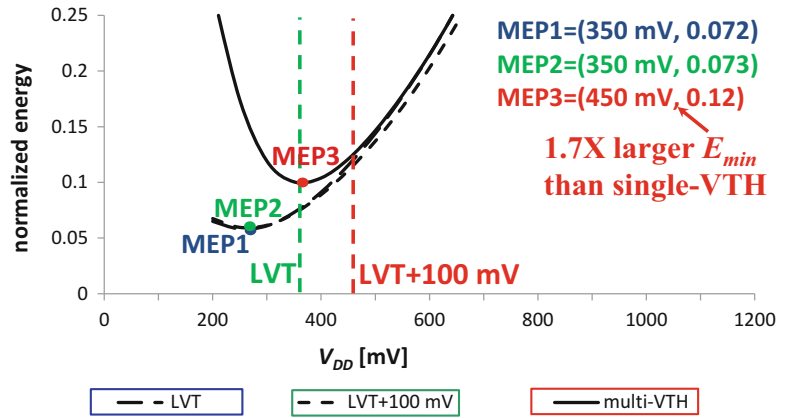
some non-trivial clocking technique to avoid timing violations under the unavoidably large variations (see Sect. 4.7), such as 2-phase latch clocking, custom latches with embedded logic, aggressive hold fix buffer insertion, and shallow clock network (3 levels, for the reasons clarified in Sect. 4.10).

The resulting energy curve versus V_{DD} for different logic depths in the reference circuit in Fig. 4.30 is reported in Fig. 4.31 in 28-nm CMOS. As expected, increasing the logic depth to the practical lower bound of 25 $FO4$ to larger logic depths of 50 and 100 $FO4$ leads to a significant 20% and a considerable 60% energy increase at the MEP. This is respectively due to a 2X and 4X leakage energy increase, due to the larger logic depth from (4.10). From (4.14a), (4.14b), such increase in E_{lkg} leads to a 2X (4X) increase in $ILDR$, which from (4.15) translates into an increase in the MEP voltage of approximately 35 mV and 65 mV (Fig. 4.31 discretizes voltages in 50-mV step, and hence results to 50 and 100 mV).

Finally, it should be observed that true minimum-energy operation actually requires a complex optimization that involves logic depth, voltage and transistor sizing. Unfortunately, no thorough methodology and no CAD support is currently available for this purpose, hence such joint optimization is still an open research question. A qualitative treatment of this problem will be presented in Sect. 4.11, to gain an insight into this fundamental design problem.

¹¹ Here, “fast” refers to the clock cycle normalized to $FO4$ (i.e., LD_{eff}), rather than the absolute clock cycle. This choice is motivated by the need for characterizing the design regardless of the specific voltage and hence $FO4$. Indeed, low values of $T_{CK}/FO4$ identify designs that would be fast at nominal and any other voltage, regardless of $FO4$. On the other hand, ultra-low voltage operation makes the absolute T_{CK} large simply because of the increase in $FO4$, not because of the design itself.

Fig. 4.32 Energy normalized to value at nominal voltage vs. V_{DD} for single- and multi-V_{TH} design (50% HVT, 50% LVT cells, logic depth of 25FO410% activity)



4.4.2 Impact of Threshold Voltage and Activity

The impact of the threshold voltage on the reference circuit in Fig. 4.30 is shown in Fig. 4.32. As expected from Sect. 4.3.4, the MEP voltage and energy are the same for different threshold voltages under a single- V_{TH} design, being in the sub-threshold region. On the other hand, mixing the two threshold voltages leads to a substantially larger MEP energy (by 1.7X in this specific case). This suggests that multi- V_{TH} design is not really advantageous at near- and sub-threshold voltages, and it should hence be avoided. Thorough analysis and justification of this observation will be provided in the next section.

The effect of activity is depicted in Fig. 4.33, which once again confirms that the MEP moves to the left when the dynamic energy is increased, as was observed in Fig. 4.21. More quantitatively, the increase in the activity factor from 3% to 10% (20%) leads to a 3.3X (6.6X) decrease in $ILDR$, which from (4.15) translates into a decrease in the MEP voltage of 51 mV and 81 mV (the latter is not precisely visualized in Fig. 4.33, as voltages are discretized in 50-mV step).

To provide a broader view on the impact of the above parameters onto the MEP position, Fig. 4.34a–c plot the statistical distribution of the MEP voltage for several different activities and logic depths, respectively for a very low, relatively low and relatively high V_{TH} . From this figure, the MEP lies in the sub-threshold

region for most of the designs, and it is pushed into in the near-threshold region only for very low threshold voltages (see Fig. 4.34a). Figure 4.34d–f show the contribution of the leakage energy as a fraction of the overall energy for the same threshold voltages. From this figure, E_{lkg} accounts for 40% of the total energy or more in most of the designs, and tends to be larger under lower threshold voltages. According to Fig. 4.23, this is because the MEP is pushed to near-threshold voltages at low V_{TH} , and the resulting E_{lkg} can be as high as 70% of the total energy in some designs (see Fig. 4.34d).

4.5 Ineffectiveness of Traditional Leakage Reduction Techniques

This section shows that traditional leakage reduction techniques (e.g., stacking, multi- V_{TH}) are far less effective at near-threshold voltages, thus posing a challenge on leakage management at such voltages.

Transistor stacking has been extensively exploited to reduce leakage in above-threshold circuits (Narendra and 2006), as the off-stacking factor¹² is typically much higher than the on-stacking factor. In other words, the series

¹²The off (on) stacking factor is defined as the factor by which the transistor current of an off (on) single transistor is reduced due to the series connection of multiple transistors having the same size.

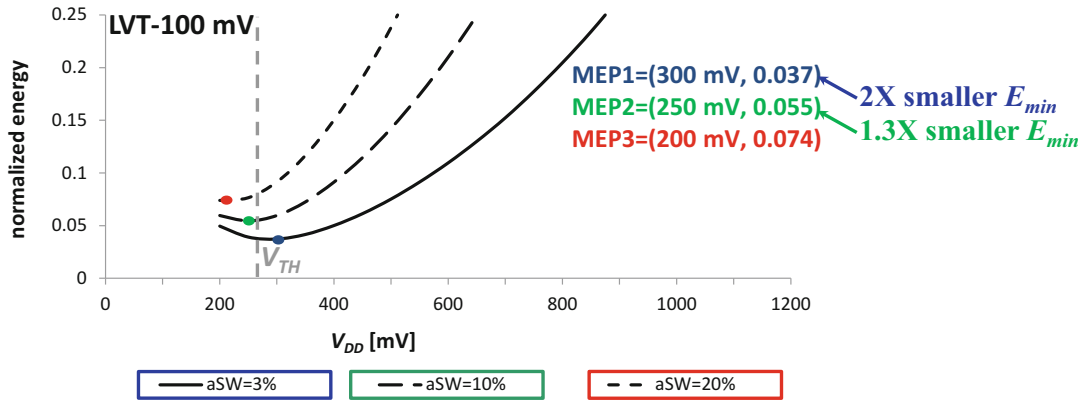


Fig. 4.33 Energy normalized to value at nominal voltage vs. V_{DD} for different activity values (logic depth of $25FO4$)

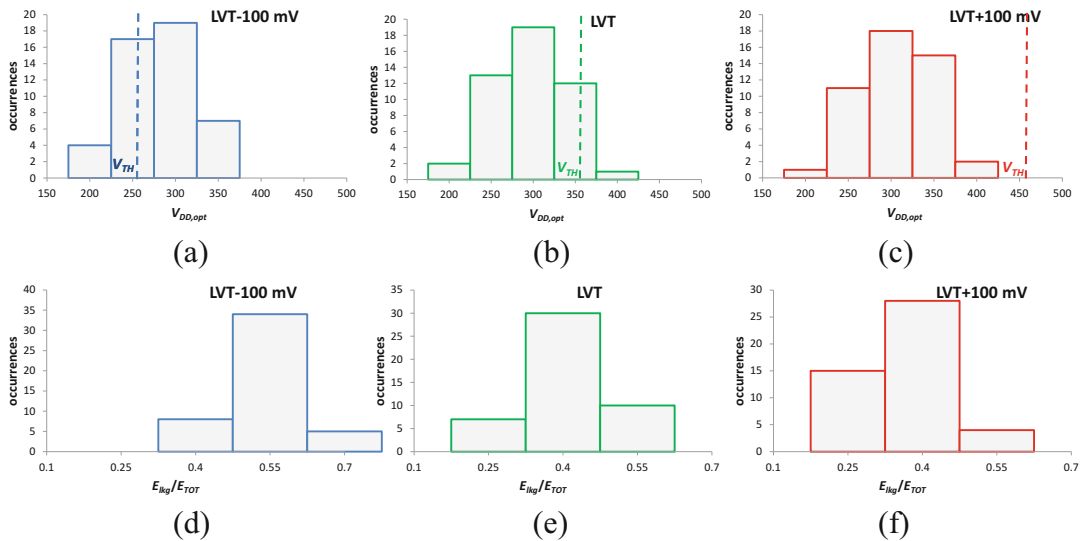


Fig. 4.34 Statistics on the MEP voltage $V_{DD,opt}$ for the reference circuit in Fig. 4.30 across different values of activity and logic depth for three threshold voltages (a)–(c). Ratio between leakage energy and total energy for same threshold voltages (d)–(f)

connection of multiple transistors reduces the leakage current much more heavily than the on-current (i.e., performance). As shown in Fig. 4.35, this is true in the above-threshold region, where the off-stacking factor for 2 to 4 transistors is larger than the on-stacking factor by an order of magnitude. At lower voltages, the on-stacking factor tends to moderately increase, for the reasons discussed in Sect. 4.2.3. At the same time, the off-stacking factor decreases exponentially when reducing V_{DD} (Narendra and Chandrakasan 2006). Indeed, the off-stacking factor for two stacked transistors can be expressed as

$e^{\alpha_{X_{off}} \cdot V_{DD}/(nkT/q)}$ (Narendra and Chandrakasan 2006), where

$$\alpha_{X_{off}} = \lambda_{DIBL} \cdot \frac{1 + \lambda_{DIBL}}{1 + 2\lambda_{DIBL}} \ll 1 \quad (4.25)$$

and approximately the same dependence is observed for a larger number of stacked transistors, as shown in Fig. 4.35. In other words, the off-stacking factor $X_{stack,off}$ is proportional to $e^{\alpha_{X_{off}} \cdot \frac{V_{DD}}{nkT/q}}$ (Narendra and Chandrakasan 2006), hence it can be expressed as

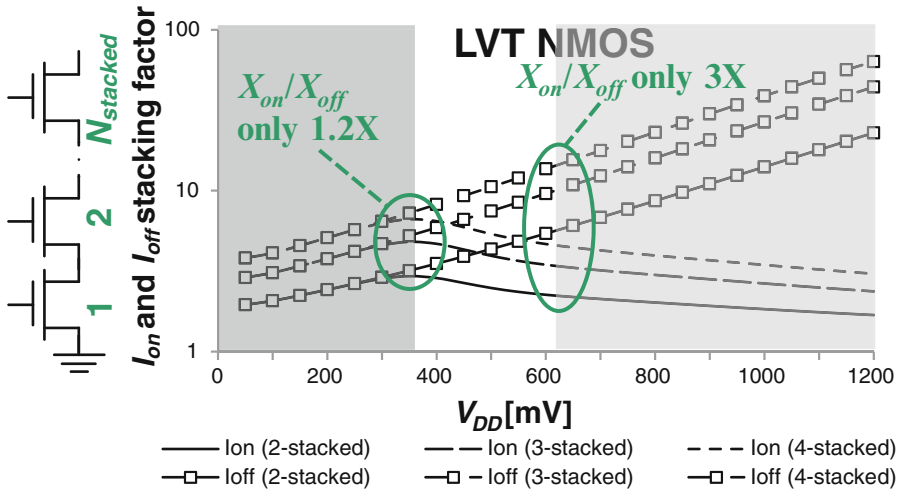
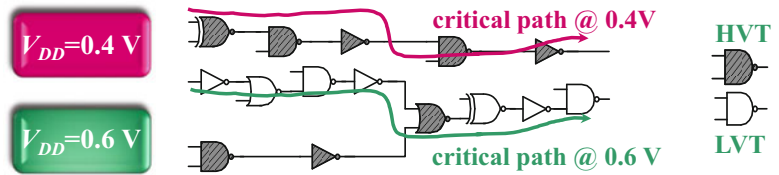


Fig. 4.35 On- and off-stacking factor X_{on} and X_{off} vs. V_{DD} for 2, 3 and 4 stacked transistors in 28 nm

Fig. 4.36 Multi- V_{TH} approach and critical path shift from LVT to HVT paths when scaling down V_{DD}



$$X_{stack,off} = X_{stack,off} \Big|_{V_{DD,nom}} \cdot e^{\alpha_{X_{off}} \frac{V_{DD} - V_{DD,nom}}{n \cdot kT/q}} \quad (4.26)$$

which tends to be very accurate across all voltages (within 2% in 28 nm, according to Fig. 4.35).

From Fig. 4.35, the off-stacking factor at near-threshold voltages is no longer much larger than the on-stacking factor, hence no significant leakage reduction is actually allowed for a given performance penalty. In other words, transistor stacking is rather ineffective in counteracting leakage at near-threshold voltages, as opposed to common low-power wisdom (i.e., above threshold).

As another traditional leakage reduction technique, let us consider the adoption of multiple threshold voltages, as depicted in Fig. 4.36 for the simple case of a design with two thresholds (i.e., low and high V_{TH}). In multi- V_{TH} designs, cells in critical paths are LVT to meet the performance requirement, whereas cells in non-critical paths are replaced by the HVT counterparts. At above-threshold voltages, such replacement does

not really degrade performance thanks to its weak dependence on V_{TH} , while it certainly reduces the leakage current thanks to its strong dependence on V_{TH} from (4.8). In other words, the multi- V_{TH} approach offers a favorable tradeoff between performance and leakage in traditional low-power designs operating above threshold. On the contrary, performance becomes very sensitive to V_{TH} at near-threshold voltages as discussed in Sect. 4.2.4, and the HVT cells are slowed down much more substantially than LVT when V_{DD} is dynamically down-scaled (see Figs. 4.5 and 4.6). As a consequence, non-critical HVT paths at a given voltage (e.g., 0.6 V in Fig. 4.36) actually become critical¹³ when down-scaling V_{DD} (e.g., 0.4 V in

¹³This is unavoidable in real designs, as overall energy-performance optimization aims to equalize the delay of different paths (De Micheli 1994), so that non-critical paths can be down-sized to reduce their energy, while maintaining the same performance target (Narayanan et al. 2010).

Fig. 4.36). In other words, the clock cycle of a multi- V_{TH} design at lower voltages is significantly larger than a single-LVT design, thus leading to a leakage energy increased compared to the latter one, from (4.10). At the same time, the leakage current of a multi- V_{TH} design is significantly larger than a single-HVT design, as the LVT cells in the design have a considerably larger leakage (typically more than an order of magnitude increase when moving from a threshold value to the immediately lower one (International Technology Roadmap for Semiconductors 2013)).

From the above considerations, multi- V_{TH} designs suffers from substantially larger leakage energy compared to single- V_{TH} designs for two concurrent reasons, under dynamic voltage scaling. Hence, multi- V_{TH} approaches actually deteriorate the energy efficiency of VLSI circuits, and should be always avoided in favor of single- V_{TH} designs. The choice of the single V_{TH} has been discussed in Sect. 4.2.4. As an example, this is shown in Fig. 4.32, where the multi- V_{TH} design of the reference circuit in Fig. 4.30 is found to be 1.7X less energy efficient than the single- V_{TH} designs. In terms of energy-performance tradeoff at the MEP, Fig. 4.37 confirms that the multi- V_{TH} design is essentially as slow as the single-HVT

design, in spite of its significantly larger energy consumption.

Similar considerations hold for other traditional leakage reduction techniques, such as power gating (Flynn et al. 2007). At above-threshold voltages, power gating is well-known to provide substantial leakage reductions due to two different mechanisms. First, the sleep transistor size can be much lower than the overall effective transistor width of the power gated circuit, as only a fraction of the cells are active at a given time. Since the relative strength of the sleep and the power gated transistors is maintained at low voltages, this reduction mechanism is essentially maintained at near-threshold voltages. Second, the sleep transistor (see Fig. 4.38a) is able to provide its large on-current during active mode ($\overline{sleep} = 0$), whereas it delivers only its off-current during sleep mode ($\overline{sleep} = 1$). Such reduction is clearly more pronounced for larger I_{on}/I_{off} ratio, which is traditionally obtained by using HVT devices for sleep transistors at above-threshold voltages. At near-threshold voltages, the transistor I_{on}/I_{off} ratio is severely degraded (by 1–2 orders of magnitude) as shown in Fig. 4.38b. Hence, the leakage reduction enabled by power gating at near-threshold voltages is worsened by at least

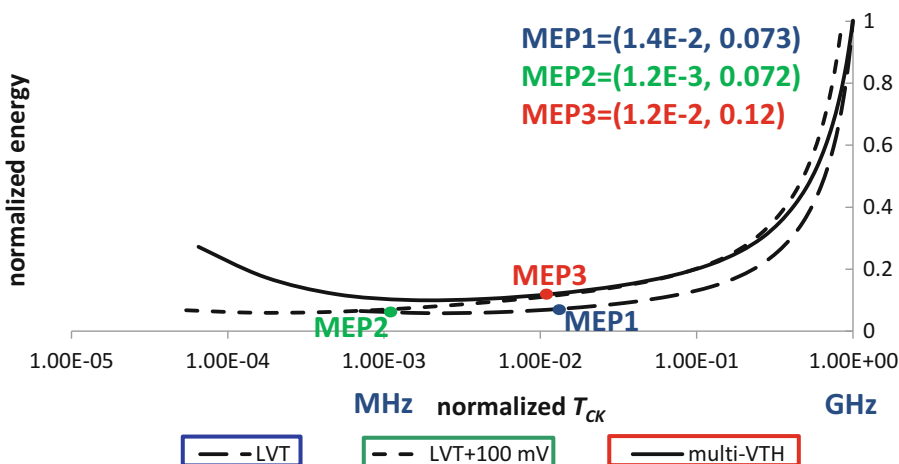


Fig. 4.37 Energy vs. clock frequency (both normalized to value at nominal voltage) for single- and multi- V_{TH} design, and logic depth of $25F04$

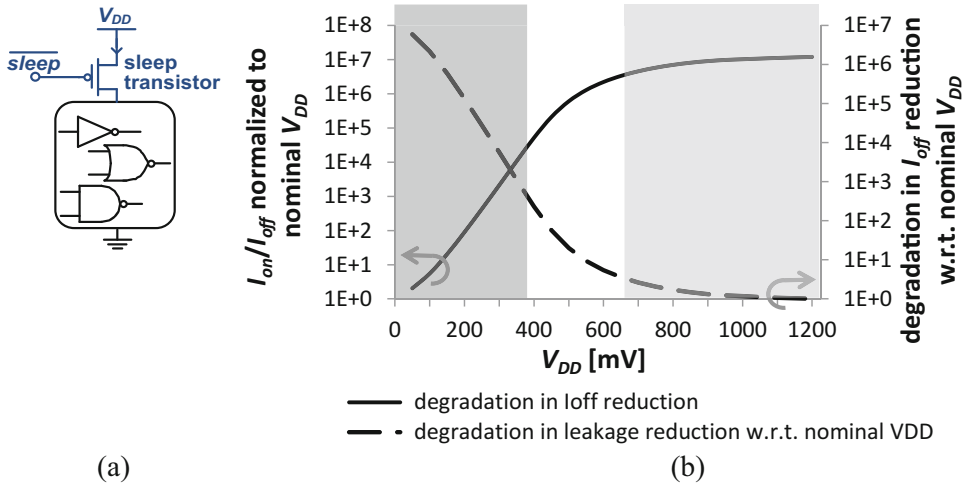


Fig. 4.38 (a) Power gating scheme, (b) I_{on}/I_{off} ratio of sleep transistor in 28 nm

one order of magnitude, compared to above threshold. Such degradation in the effectiveness of power gating at near-threshold voltages can be partially recovered by boosting the gate voltage of the sleep transistor (Myers et al. 2015). Indeed, boosting its gate voltage only during active mode significantly increases I_{on} , while maintaining the same I_{off} . At near-threshold voltages, the sleep transistor I_{on}/I_{off} ratio (and hence the effectiveness of power gating) can be further improved by using thick-oxide (i.e., I/O) NMOS transistors whose gate is powered at the large I/O voltage (e.g., 1.8 V instead of 1 V). In this case, such I_{on}/I_{off} improvement is achieved at the expense of a larger energy and slower transient to turn on the sleep transistor, and hence to switch from sleep to active mode.

4.6 Challenges: Performance

As discussed in Sect. 4.2, operation at near-threshold voltages entails a $\sim 10X$ penalty in terms of $FO4$ and hence performance, compared to the same architecture operating at nominal voltage. For sub-100 nm technologies, $FO4$ at near-threshold voltages is typically in the order of few hundreds of picoseconds. For reasonable architectures with a logic depth of up to several tens of $FO4$, this translates into a cycle time in

the order of nanoseconds. Hence, throughputs in the order of hundreds of MOPS (Millions of Operations per Second) are easily achievable by near-threshold microprocessor cores. Such level of performance achievable near the threshold is actually acceptable for (or can exceed) the typical requirements of IoT systems, at least in the most frequent operation modes and in most of the practical applications. Higher performance might be needed occasionally in some applications, or customarily for compute-intensive ones, such as computer vision or real-time pattern recognition.

Sustained throughputs that are higher than hundreds of MOPS can always be obtained through appropriate architectures at near-threshold voltages (e.g., multi-core) and specialized hardware, as discussed in Chap. 3. Occasional performance boosts can be achieved through wide dynamic voltage scaling, i.e., by raising V_{DD} from the MEP to the nominal voltage (Chandrakasan et al. 2010). Such temporary voltage up-scaling permits to increase the performance by one (two) order(s) of magnitude, when the MEP is in the near-threshold (sub-threshold) region. This performance increase is achieved at the expense of an increase in the energy per operation, as summarized in Fig. 4.39 for several integrated prototypes (Abouzeid et al. 2012; Gammie et al. 2011; Hsu et al. 2012; Jain et al. 2012; Kaul et al. 2012;

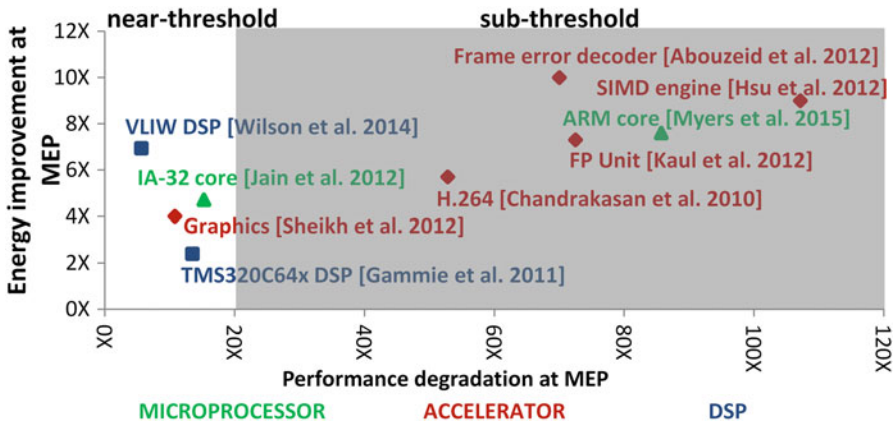
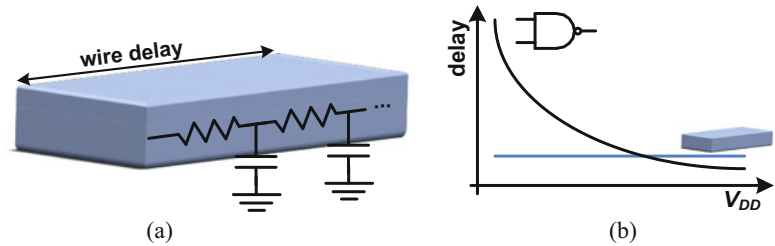


Fig. 4.39 Energy improvement at MEP vs. performance degradation at MEP, as compared to operation at nominal voltage

Fig. 4.40 (a) RC wire delay, (b) relative scaling of gate and wire delay vs. V_{DD}



Myers et al. 2015; Sheikh et al. 2012; Wilson et al. 2014; Jacquet et al. 2013). From this figure, this energy increase is more pronounced when the MEP is in the sub-threshold region, due to the larger voltage difference between the MEP and the nominal voltage.

As another property of circuits operating near the MEP, the gate delay dominates over the wire delay, as shown in Fig. 4.40. Indeed, they might be comparable at nominal voltage in realistic VLSI architectures, due to the significant resistive, capacitive, and sometimes inductive parasitics of metal wires. However, operation at the MEP voltage determines a substantial increase in the gate delay, while keeping the wire delay constant. Hence, the wire delay is no longer a challenge in circuits with nearly-minimum energy, which certainly simplifies the design, the circuit modeling and the timing closure.

The reduced wire-to-gate delay ratio around the MEP has also important consequences on the

choice of the architecture, and the way the latter is mapped into the physical level. Indeed, VLSI architectures for nearly-minimum energy need to be different from traditional low-power architectures (i.e., for above-threshold operation). More specifically, signals can be propagated through a wider silicon area compared to nominal voltage operation. In detail, for unrepeated wires a $\sim 3X$ longer distance¹⁴ can be covered by the same wire at near-threshold voltages, as compared to the same circuit operating at nominal voltage, when maintaining the wire delay a fixed fraction of the clock cycle. For similar reasons, repeated wires require 3X fewer repeaters per wire unit length since the optimal distance between

¹⁴This is due to the well-known quadratic dependence of the RC wire delay on its length (Weste and Harris 2011), and assuming a 10X $FO4$ degradation at the MEP compared to the nominal voltage.

repeaters is proportional to the square root of $FO4$ (Weste and Harris 2011), thus improving the route-ability. At the same time, the size of each repeater at MEP needs to be increased by 3X compared to nominal voltage, as its performance-optimal size is proportional to $FO4$ (Weste and Harris 2011). Since the number of repeaters is reduced by the same factor by which their area is increased, the energy and area cost of intra-chip global communication in designs around the MEP remains approximately the same. In summary, VLSI architectures for nearly-minimum energy can afford more global communications and larger modules (e.g., shared caches), compared to traditional low-power architectures. Such profound difference in the communication-computation energy/performance tradeoff requires the adoption of innovative architectures, as discussed in Chap. 3.

around the MEP. In general, process, voltage and temperature variations as well as aging impose an additional timing margin that stretches the clock cycle as shown in Fig. 4.41. This conservative approach preserves correct functionality and performance specifications even in the worst-case die and environmental conditions.

The above cycle time margin resulting from variations translates into an increase in the energy per operation, as faster chips are forced to operate as slowly as the worst-case die and at the same voltage (which is higher than needed). Figure 4.42 shows the voltage increase required by the circuit in Fig. 4.30 to maintain a given clock cycle under a given clock cycle margin, as well as the resulting energy increase, assuming a logic depth of $25FO4$, 10% activity, and LVT transistors. From Fig. 4.42, the voltage increase imposed by variations is fairly linear with the cycle time margin, and tends to be larger at higher nominal operating voltages. This is because $FO4$ (i.e., the cycle time from Sect. 4.3.2) is less sensitive to voltage increases at higher operating voltages, and hence requires larger increase to achieve a given percentage

4.7 Challenges: Variations

In this section, the impact of variations is analyzed in the context of circuits operating

Fig. 4.41 Nominal cycle time and additional margin accounting for variations

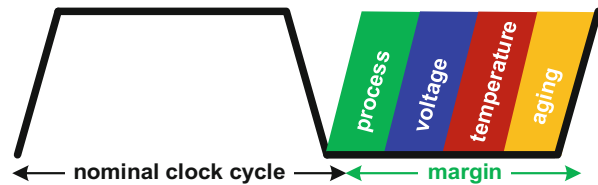
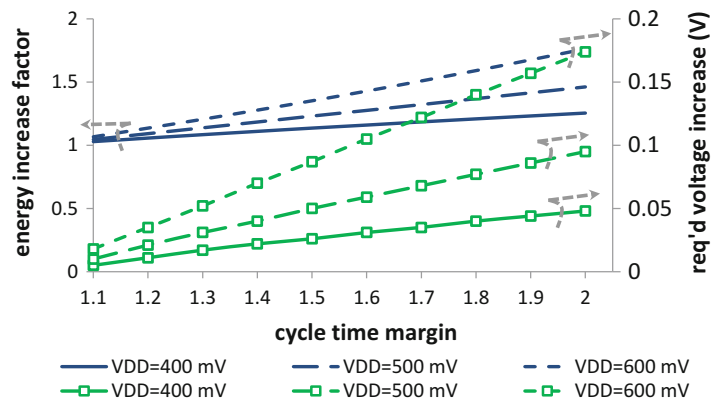


Fig. 4.42 Required V_{DD} to sustain a given performance specification vs. cycle time margin (i.e., factor by which the cycle time needs to be increased due to variations), and resulting energy increase due to variations



performance improvement. From Fig. 4.42, a typical 1.5X–2X cycle time increase requires a voltage increase by 100–200 mV to sustain the same performance as the nominal corner, which leads to an energy increase by a factor of 1.5X–2X. In other words, variations at near threshold entail a very large energy cost, which can negate the advantages of operating at the MEP. Accordingly, variations need to be accounted for in first place in the design of near-threshold circuits rather than an afterthought, as discussed more in detail in the following. Similar considerations hold for the sensitivity to soft errors, which is somewhat increased at near threshold voltages, compared to nominal voltage. On the other hand, operation at near-threshold voltages suppresses most of aging and reliability issues, such as Bias Temperature Instability, Hot Carrier Injection, Time Dependent Dielectric Breakdown. Indeed, such phenomena are all exponentially dependent on the supply voltage, and its reduction to near-threshold voltage substantially mitigates them.

4.7.1 Process Variations

Random (within-die) process variations are well known to be responsible for a major fraction of the cycle time margin, having a much heavier effect than fully correlated (die-to-die) variations (Orshansky et al. 2008). Indeed, threshold voltage variations at low voltages are dominated by random dopant fluctuations (Alioto et al. 2010; Orshansky et al. 2008), and their effect requires much more sophisticated feedback schemes that are immune to transistor mismatch (e.g., with

timing error detection or prediction (Bowman et al. 2009; Bowman et al. 2011; Das et al. 2009; Khayatzadeh et al. 2016; Zhang et al. 2016)), rather than corner-based adaptive voltage scaling and body biasing techniques (Gregg and Chen 2007; Meijer and Pineda de Gyvez 2012; Martin et al. 2002; Olivieri et al. 2005; Tschanz et al. 2002). Accordingly, our analysis in the following will be focused on random variations.

At low voltages, process variations determine a much larger path delay variations than above-threshold voltages due to two phenomena:

1. the variability of the gate delay defined as the ratio between the standard deviation and mean value increases significantly
2. the probability distribution function (PDF) of such delay is no longer Gaussian for short paths, and has a longer tail on the right side.

Regarding the first phenomenon, the variability of the critical path delay is mainly due to the intrinsically larger variability of the transistor I_{on} current (see Eq. (4.6)). This is mostly due to the larger impact of the threshold voltage, and hence of its variations, at lower voltages (see Sect. 4.2.1). More quantitatively, the delay variability is approximately equal to the variability in I_{on} from (4.6). If the nominal threshold voltage V_{TH} is subject to a variation ΔV_{TH} that is Gaussian distributed with zero mean and standard deviation $\sigma_{V_{TH}}$, from (4.3) to (4.4) the variability of I_{on} for above- and sub-threshold voltages is readily found to be

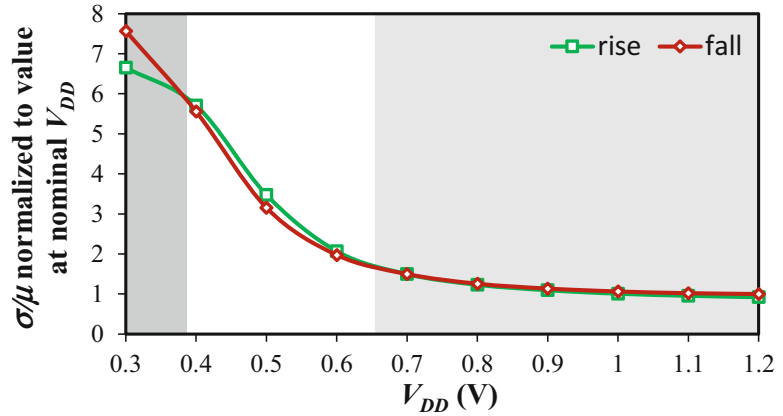
$$\left. \frac{\sigma_{\tau_{PD}}}{\mu_{\tau_{PD}}} \right|_{above-threshold} \approx \left. \frac{\sigma_{I_{on}}}{\mu_{I_{on}}} \right|_{above-threshold} = \frac{\sigma_{V_{TH}}}{V_{DD} - V_{TH}} \quad (4.27a)$$

$$\left. \frac{\sigma_{\tau_{PD}}}{\mu_{\tau_{PD}}} \right|_{sub-threshold} \approx \left. \frac{\sigma_{I_{on}}}{\mu_{I_{on}}} \right|_{sub-threshold} = \sqrt{e^{\left(\frac{\sigma_{V_{TH}}}{nV_T}\right)^2} - 1} \quad (4.27b)$$

For example, for the typical values $\sigma_{V_{TH}} = 35$ mV and $V_{TH} = 0.4$ V in 28 nm CMOS, the gate delay variability turns out to be 7% at 0.9 V, and

124% in the sub-threshold region (e.g., 0.3 V). In other words, the delay variability in sub-threshold is an order of magnitude larger

Fig. 4.43 Variability of $FO4$ normalized to value at 1.2 V vs. V_{DD} (28 nm CMOS)



than above threshold. At near-threshold voltages, the variability is somewhat intermediate.

Figure 4.43 plots the variability of the $FO4$ delay normalized to the value at nominal voltage in 28 nm CMOS. From this figure, the gate delay variability increases when V_{DD} is reduced, and becomes 2X–6X larger than at nominal V_{DD} at near threshold, and about an order of magnitude larger in the deep sub-threshold region. A typical delay variability of around 6–8% at nominal voltage in 28 nm translates into a sizable delay variability of various tens of percentage points at near threshold. To achieve a parametric yield of approximately 99%, three standard deviations are needed, hence the margin for a single gate can easily be 100%, which entails an unfeasibly large margin in Fig. 4.41 (i.e., an unacceptably high energy cost from Fig. 4.42, which easily offsets the energy benefit of operating at near-threshold voltages). When the MEP is in sub-threshold region, such margin becomes even higher.

Regarding the second phenomenon that was observed above, the statistical delay distribution of a single gate is no longer Gaussian when operating at near- and sub-threshold voltages (Alioto 2012; Alioto et al. (in press); Gammie et al. 2011). In the sub-threshold region, I_{on} and hence the gate delay are lognormally distributed due to the exponential dependence of I_{on} on the threshold voltage in (4.4), being the latter Gaussian distributed. As shown in Fig. 4.44, the

lognormal distribution has a much longer tail compared to the Gaussian distribution, at same standard deviation. This leads to a considerable increase in the number of standard deviations needed as design margin to meet a given yield target, as shown in Table 4.7. For example, from this table the worst-case gate delay margin across 99.9% of the cases is three standard deviations for Gaussian (i.e., above threshold), and twenty standard deviations for lognormal (i.e., sub-threshold). In the near-threshold region, the distribution is somewhat intermediate between above- and sub-threshold, and hence it is neither perfectly Gaussian nor lognormal. This is shown in Fig. 4.45a–c, which show the quantile-quantile (Q–Q) plot (Walpole et al. 2006) of the statistical $FO4$ delay sample in 28 nm CMOS versus the theoretical quantiles of a Gaussian distribution with same mean and standard deviation. The deviation from a straight line (i.e., perfect Gaussian behavior) of the Q–Q plot becomes noticeable at near threshold (see Fig. 4.45b), and is substantial at sub-threshold voltages (see Fig. 4.45c). Figure 4.45d confirms the $FO4$ lognormal distribution in sub-threshold.

The above considerations of non-Gaussian delay distribution at low voltages hold for single logic gates, and can be extended to short paths, i.e., paths that can be problematic in terms of hold time violations rather than setup time. Accordingly, short paths and hold fix at sub- and near-threshold voltages requires a much wider design

Fig. 4.44 Probability density function of Gaussian and lognormal distribution at same standard deviation ($\sigma = 1$)

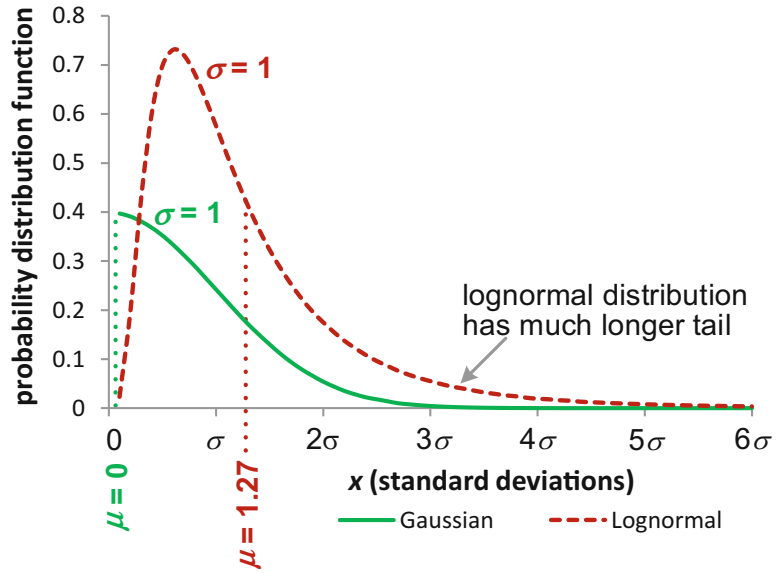


Table 4.7 Number of V_{TH} Standard Deviations beyond the Mean to Achieve Given Yield Target in Gaussian and Lognormal

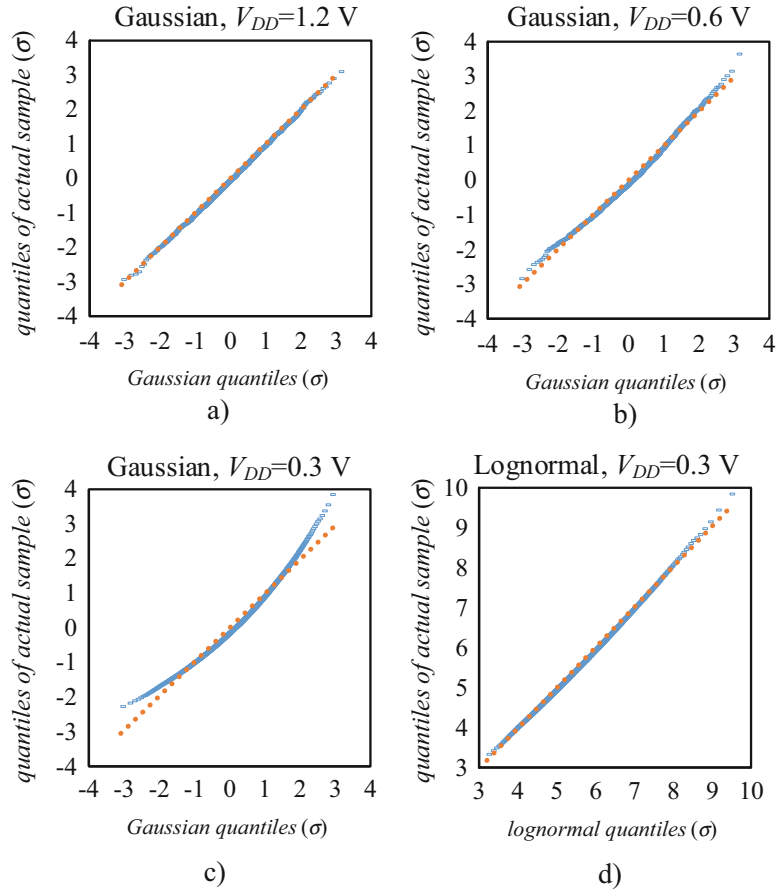
Yield target	Single gate		Logic path (lognormal gate delay)			
	Gaussian	Lognormal	Logic depth ↓			
			4 gates	8 gates	16 gates	32 gates
84%	$\sigma_{V_{TH}}$	$e \cdot \sigma_{V_{TH}}$	$1.07 \cdot \sigma_{V_{TH}}$	$1.05 \cdot \sigma_{V_{TH}}$	$1.03 \cdot \sigma_{V_{TH}}$	$1.02 \cdot \sigma_{V_{TH}}$
97.7%	$2\sigma_{V_{TH}}$	$e^2 \cdot \sigma_{V_{TH}} \approx 7.4 \cdot \sigma_{V_{TH}}$	$2.07 \cdot \sigma_{V_{TH}}$	$2.05 \cdot \sigma_{V_{TH}}$	$2.03 \cdot \sigma_{V_{TH}}$	$2.02 \cdot \sigma_{V_{TH}}$
99.87%	$3\sigma_{V_{TH}}$	$e^3 \cdot \sigma_{V_{TH}} \approx 20.1 \cdot \sigma_{V_{TH}}$	$3.07 \cdot \sigma_{V_{TH}}$	$3.05 \cdot \sigma_{V_{TH}}$	$3.03 \cdot \sigma_{V_{TH}}$	$3.02 \cdot \sigma_{V_{TH}}$
99.997%	$4\sigma_{V_{TH}}$	$e^4 \cdot \sigma_{V_{TH}} \approx 54.6 \cdot \sigma_{V_{TH}}$	$4.07 \cdot \sigma_{V_{TH}}$	$4.05 \cdot \sigma_{V_{TH}}$	$4.03 \cdot \sigma_{V_{TH}}$	$4.02 \cdot \sigma_{V_{TH}}$
99.99997%	$5\sigma_{V_{TH}}$	$e^5 \cdot \sigma_{V_{TH}} \approx 148.4 \cdot \sigma_{V_{TH}}$	$5.07 \cdot \sigma_{V_{TH}}$	$5.05 \cdot \sigma_{V_{TH}}$	$5.03 \cdot \sigma_{V_{TH}}$	$5.02 \cdot \sigma_{V_{TH}}$
99.9999999%	$6\sigma_{V_{TH}}$	$e^6 \cdot \sigma_{V_{TH}} \approx 403.4 \cdot \sigma_{V_{TH}}$	$6.07 \cdot \sigma_{V_{TH}}$	$6.05 \cdot \sigma_{V_{TH}}$	$6.03 \cdot \sigma_{V_{TH}}$	$6.02 \cdot \sigma_{V_{TH}}$

margin, compared to above-threshold. In other words, the timing margin against hold time violations at low voltages tends to be very large compared to nominal voltage, and hence requires a much larger number of hold fix buffers.

On the other hand, long logic paths have a Gaussian delay distribution even in sub-threshold voltages. This is because of the Central Limit theorem, which guarantees that the sum of non-Gaussian random variables rapidly tends to a Gaussian distribution, when increasing the number of variables being summed (Walpole et al. 2006) (e.g., the number of logic gates whose delays are added to derive the critical path delay). This is quantitatively

shown in Table 4.7 for 4, 8, 16 and 32 equal cascaded gates, which are individually assumed to have a lognormal delay distribution, as relevant to the sub-threshold region. Indeed, this table shows that margin in terms of standard deviations is essentially the same as an ideal Gaussian distribution even for a relatively short path of 4 cascaded gates, and is closer for a larger number of gates. This means that the clock cycle distribution is Gaussian at any voltage, and hence the margin in terms of standard deviations is the same as nominal voltage. In other words, the timing margin against setup time violations at low voltages scales like $FO4$, as opposed to hold violations.

Fig. 4.45 Q–Q plots of $FO4$ distribution (y axis) in 28 nm CMOS at (a) 1.2 V, (b) 0.6 V, (c) 0.3 V with normal distribution on the x axis. (d) at 0.3 V with lognormal distribution on the x axis (100,000 Monte Carlo runs)



4.7.2 Voltage and Temperature Variations

Voltage variations have a heavy impact on the clock cycle margin, due to the strong sensitivity of I_{on} and hence the gate delay on V_{DD} (see Sect. 4.2.1). Figure 4.46 plots the cycle time margin associated with a typical 5% and 10% voltage drop of the circuit in Fig. 4.30 in 28 nm CMOS. As all gate delays scale approximately like $FO4$ when V_{DD} changes, this example is representative of any logic path. From this figure, a large cycle time margin of 20–50% is imposed by supply variations, if not kept under strict control. As discussed in Sect. 4.10, supply variations in systems designed for nearly-minimum energy operation are dominated by fluctuations in the output voltage of the regulator

providing the supply. Accordingly, supplies for minimum-energy operation need to be designed with quite stringent specifications on voltage stability across temperatures, as well as line and load regulation.

Temperature variations in circuits designed for minimum energy have an effect that is quite different from above-threshold circuits. Indeed, larger temperatures lead to a substantial increase in the energy per operation at low voltages, due to the large contribution of the leakage energy (see Sect. 4.3.2). Such effect is more pronounced in architectures with larger leakage energy, e.g., with larger logic depth. Figure 4.47 plots the energy versus V_{DD} of the circuit in Fig. 4.30 at different temperatures (27 °C and 70 °C), and for logic depths widely ranging from 25 $FO4$ to 100 $FO4$. From this figure, the minimum energy

Fig. 4.46 Plot of the cycle time margin vs. V_{DD} (assuming cycle time scaling to be proportional to $FO4$)

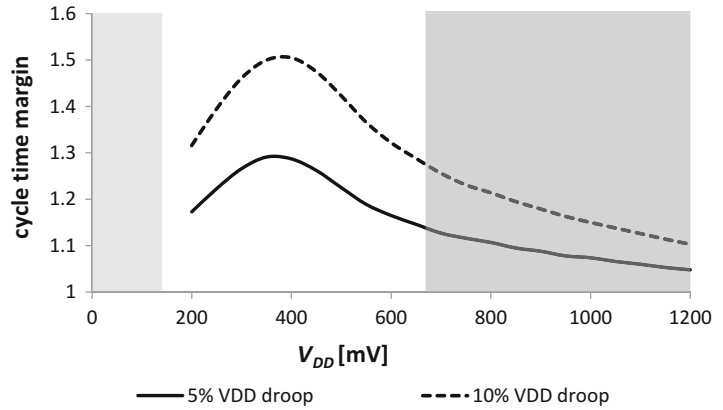
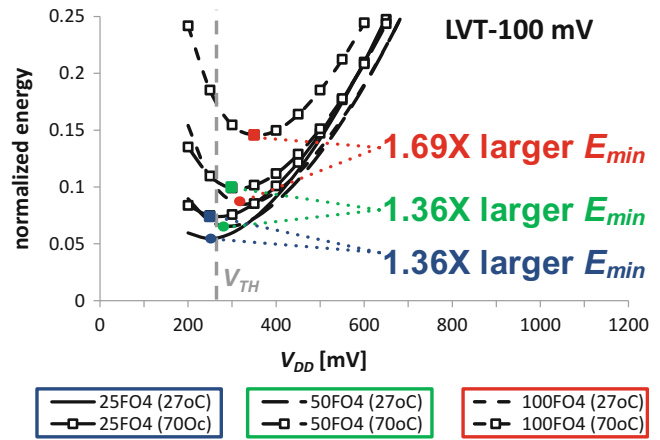


Fig. 4.47 Impact of temperature on minimum energy vs. V_{DD} for different logic depths (10% activity, very low V_{TH})



is heavily influenced by the operating temperature, as it is increased by a factor of 1.3X for well-designed architectures with reasonable logic depth, and 1.7X for less energy-efficient and leakier architectures.

As further difference compared to traditional above-threshold low-power designs, the performance of circuits operating around the MEP actually benefits from increased temperature. Indeed, the I_{on} transistor current is much more sensitive to the threshold voltage rather than the carrier mobility, hence it increases at larger temperatures. Figure 4.48 shows the $FO4$ delay versus V_{DD} for various threshold voltages. From this figure, a temperature raise from 27 °C to 70 °C leads to a 1.4X–2X $FO4$ reduction at V_{DD} equal to the threshold voltage V_{TH0} in

(4.7). Such effect is less pronounced at higher threshold voltages, as the corresponding higher supply voltage emphasizes the carrier velocity saturation and the mobility degradation due to high-field operation, which in turn weaken the dependence of I_{on} on V_{TH0} . At above-threshold voltages around 700–800 mV, the effect of temperature is insignificant. At larger voltages, the temperature has the traditional inverse effect on the performance, and is much weaker (e.g., 2% change due to a temperature change from 27 °C to 70 °C) than near threshold. Hence, unless the operating temperature range set by the application is narrow (e.g., indoor applications), active compensation of temperature variations is essential in any integrated system aiming at minimum-energy operation.

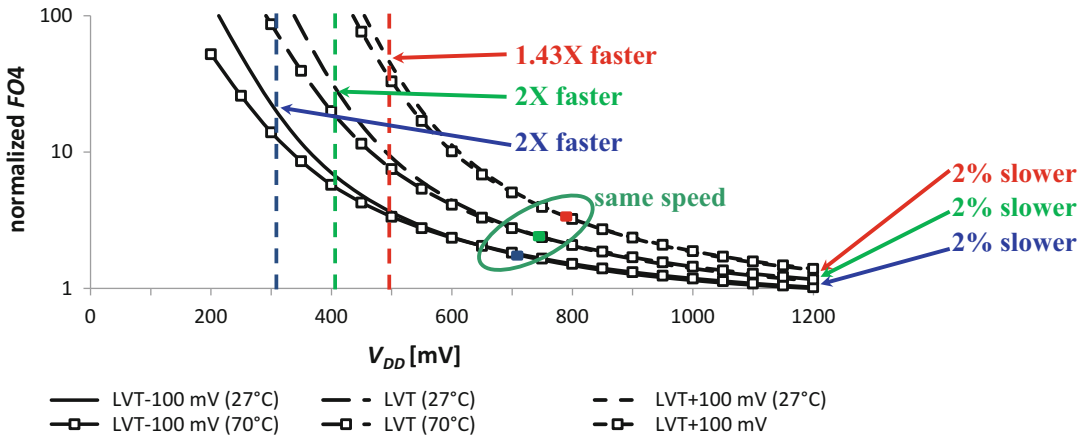


Fig. 4.48 Impact of temperature on performance vs. V_{DD} for different threshold voltages ($FO4$ is normalized to the value at nominal voltage for lowest V_{TH})

Table 4.8 Dependence of the delay variability in logic paths, cells and transistors

		Design parameter	Delay variability dependence
Logic path		logic depth LD	$\propto \frac{1}{\sqrt{LD}}$
Standard cell		$N_{stacked}$	$\propto \frac{1}{\sqrt{N_{stacked}}}$
Transistor		channel width W	$\propto \frac{1}{\sqrt{strength}}$

4.8 The Leakage-Variability Tradeoff

Operation around the MEP introduces a tradeoff that is not encountered in traditional low-power above-threshold designs, namely the variability-leakage tradeoff. This is an unavoidable tradeoff that constrains the design at all levels of abstraction and is tightly linked to the averaging effect of additive variations, as discussed below.

At the gate level, a logic path with logic depth LD has a delay that is the sum of LD delays, as depicted in Table 4.8. The resulting path delay variability is inversely proportional to \sqrt{LD} thanks to the averaging effect of the random variations across cascaded gates (Alioto et al. 2010; Merrett et al. 2010) (i.e., more cascaded gates reduce the overall delay variability thanks

to better averaging across a larger number of cells). Hence, the reduction of the delay variability would require the adoption of microarchitectures with larger logic depths. On the other hand, larger logic depths increase the clock cycle and hence the leakage energy per cycle from (4.10). In other words, the mitigation of delay variations comes at the cost of a higher leakage energy, and vice versa. Such tradeoff is very specific to operation at near- and sub-threshold voltages, due to the much more important contribution of the leakage energy, as opposed to above-threshold designs.

At the cell circuit level, a similar tradeoff is encountered when the number of stacked transistors is considered in a standard cell (Alioto et al. 2010; Merrett et al. 2010) (i.e., the cell fan-in). Indeed, the variability of the I_{on} current delivered by the cell to the load, and hence the

cell delay variability, is inversely proportional to $\sqrt{N_{stacked}}$ as in Table 4.8. Thus, the mitigation of variations through more stacked transistors comes at the cost of larger delay (see considerations on stacking in Sect. 4.2.3) and hence larger leakage energy from (4.10).

At the transistor level, from Table 4.8 wider transistors exhibit smaller I_{on} variability thanks to the Pelgrom's law (Pelgrom et al. 1989). Hence cells with larger strength have smaller delay variability, as the latter is inversely proportional to $\sqrt{strength}$. Again, the delay variability mitigation comes at the cost of larger leakage energy, as the leakage current drawn by a cell is proportional to its strength.

Summarizing, the variability-leakage tradeoff is an unescapable challenge in the design of circuits and systems operating around the MEP, as opposed to traditional low-power above-threshold designs. This tradeoff involves all levels of abstraction, and needs to be constantly taken care of during the design process. Such tradeoff can be broken by introducing innovative design techniques that do not purely rely on timing margining, as discussed later in this chapter.

4.9 Near-Threshold Cell Libraries

Designing cell libraries for operation around the MEP certainly helps manage the peculiar

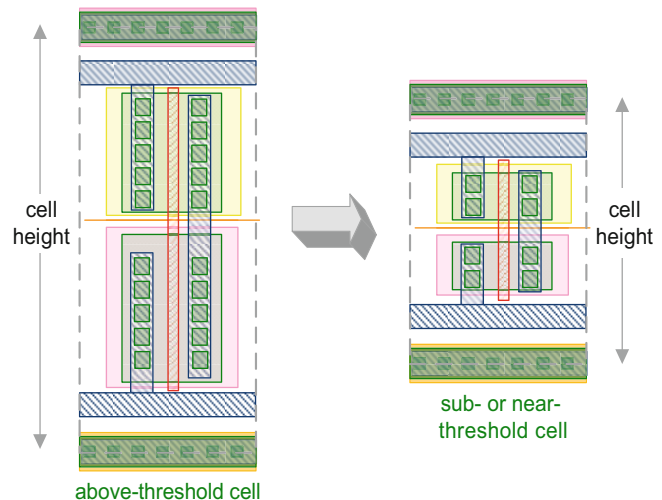
tradeoffs observed at near-threshold voltages in a more efficient manner, at the cost of additional design and characterization effort.

Since performance is not the main objective, near-threshold cell libraries can be designed with short standard cells (e.g., 7 metal tracks), as transistors typically do not need to be wide, as depicted in Fig. 4.49. In near-threshold designs, taller cells (e.g., 10–12 metal tracks) can achieve higher performance but lead to a significant area efficiency degradation, and longer interconnects, thus degrading energy efficiency.

The composition of near-threshold libraries does not need to be as wide as libraries for above-threshold (i.e., higher performance) operation. Indeed, cell versions with very large strength can be suppressed, as they are typically not used due to the more relaxed performance constraints. Similarly, cells with large fan-in need to be eliminated, since they suffer from disproportionately larger delay, as discussed in Sect. 4.2.3. Typically, libraries with around 100 cells are adequate for near-threshold designs. From observations of prior designs, the energy reduction obtained through a custom near-threshold library can be in the order of 20% (Gemmeke et al. 2013; Gammie et al. 2011), compared to a pruned out conventional library for above-threshold voltages (see below).

The circuit design of cells is affected by near-threshold operation in terms of sizing as well.

Fig. 4.49 Near-threshold cells are shorter than typical above-threshold cells



Indeed, minimum transistor size needs to be skipped in technologies that are significantly affected by Narrow Channel Effects (see Sect. 4.2.2), to avoid the related increase in the transistor threshold voltage.

Near-threshold libraries might need to be enriched with cells that are normally not available in above-threshold libraries. For example, cells with thick-oxide transistors might be needed for always-on blocks (see Chap. 1) that need to be very low leakage, or connected directly to 3.6-V LiIon batteries (see Chap. 15). Being particularly critical in terms of the minimum voltage V_{min} assuring correct operation, flip-flops usually need to be thoroughly redesigned to achieve adequate functional yield at low voltages. This is usually achieved through circuit techniques that eliminate the potential current contention between transistors (Jain et al. 2012; Kim et al. 2014a). V_{min} is further reduced by replacing conventional dynamic circuits (e.g., periphery in register files) by their static CMOS counterparts. As summarized in Fig. 4.50, V_{min} is determined by several contributions arising at the process and circuit level, and is certainly dominated by variations (Alioto 2012).

As an alternative option, existing cell libraries designed for above-threshold regions can be reused at lower voltages, after proper pruning to eliminate the cells that suffer from robustness issues or particularly pronounced delay increase (Alioto 2010; Wang et al. 2006). In a given library designed for above-threshold voltages, the number of usable cells at lower voltages decreases when reducing V_{DD} , as fewer cells can operate reliably at lower voltages. Typically, as summarized in Fig. 4.51, the suppression of cells with a high fan-in (e.g., 4) leads to approximately 100-mV V_{min} reduction (Gemmeke et al. 2013).

4.10 Clock and Supply Networks for Near-Threshold Operation

The design of clock networks for near-threshold designs is very different from above-threshold networks, due to the very different balance between clock repeater and wire delay, and the clock skew is determined by different dominant mechanisms (Alioto 2014; Lin et al. 2017; Seok et al. 2011; Tolbert et al. 2011). At above-threshold voltages, several levels of clock

Fig. 4.50 Breakdown of minimum supply voltage V_{min} of logic gates ensuring correct operation (Alioto 2012)

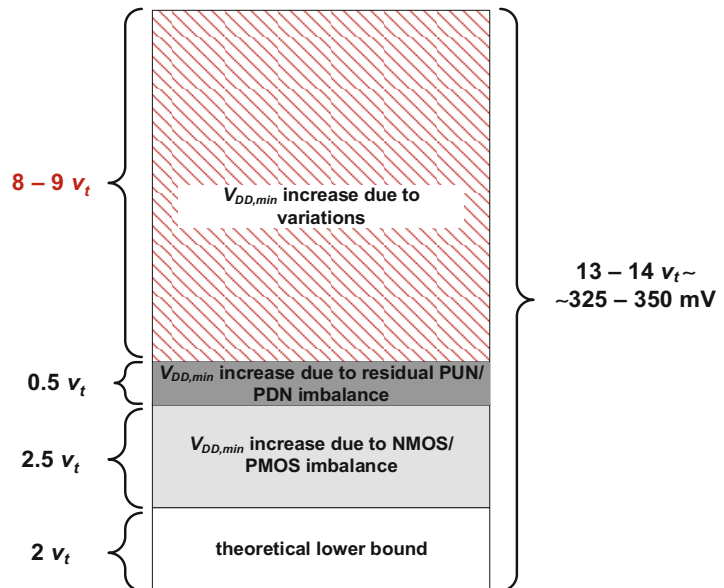


Fig. 4.51 Percentage of library cells operating correctly vs. V_{DD} (Gemmeke et al. 2013)

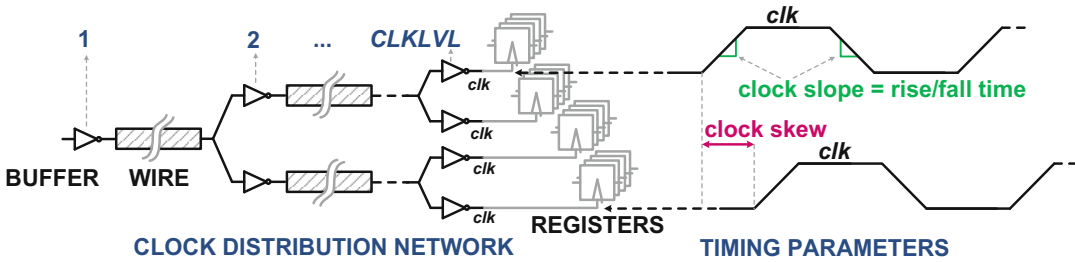
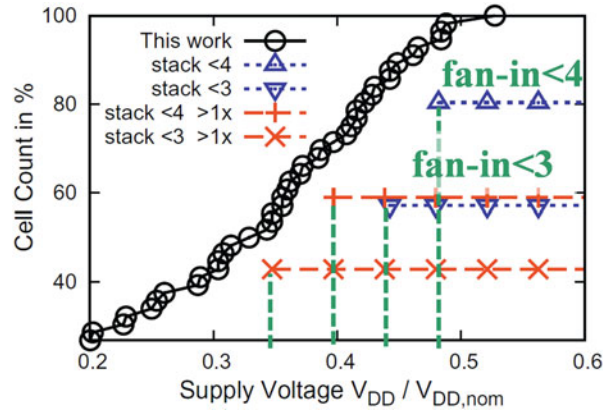


Fig. 4.52 General clock network structure and related timing parameters (Alioto 2014)

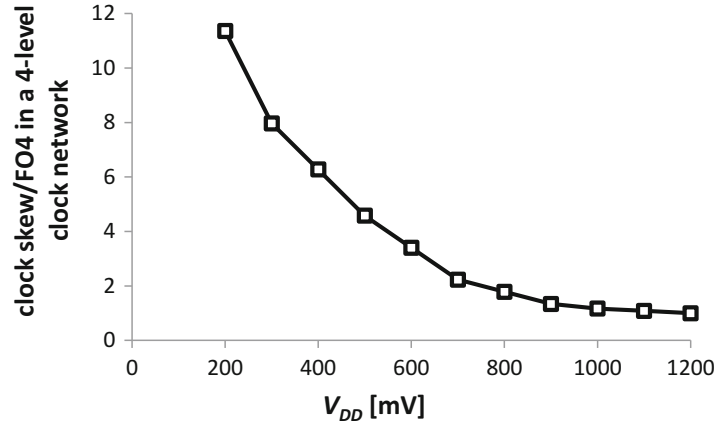
repeaters are needed to frequently interrupt wires to limit the related RC time constant and hence clock slope through the wires (Xanthopoulos 2009) (see Fig. 4.52). Indeed, excessive clock slope induces large random delay variations in the clock repeaters at intermediate nodes of the clock network, and degrades flip-flop nominal timing parameters (as well as its variations) when considering the sinks of the same network. In other words, the significant wire RC delay and its impact on clock skew through the clock slope justifies the adoption of deep above-threshold clock networks.

At sub- and near-threshold voltages, the gate delay becomes much larger than the wire delay (see Sects. 4.2.3 and 4.6), hence the clock slope through wires is no longer an issue, and the random skew is dominated by the intrinsic variations in the clock repeaters. According to the Central Limit theorem (Walpole et al. 2006), the random skew standard deviation is proportional to the square root of the number of

cascaded repeaters, i.e., of the depth of the clock network. Accordingly, shallow networks need to be used at sub- and near-threshold voltages, so that the dominant skew contribution due to the number of clock repeaters is reduced.

From the above considerations, the design the clock network at a given voltage leads to a skew degradation at the other end of the voltage range. As an example, Fig. 4.53 plots the skew of a clock network in 28 nm that has been designed at 1.2 V and used at lower voltages. This figure shows that the skew at low voltages becomes several $FO4$ and even exceeds $10FO4$ in sub-threshold. This means that the skew in a clock network used in a wide range of voltages becomes a large fraction of typical cycle time targets of energy-efficient designs (see Sect. 4.4.1). In other words, using a clock network in a wide voltage range leads to significant performance degradation (or energy efficiency, if V_{DD} is increased to recover the lost performance). Similarly, the clock skew easily exceeds

Fig. 4.53 Clock skew of a sample clock network in 28 nm designed at 1.2 V (normalized to $FO4$) (Alioto 2014)



the available hold margin (Alioto et al. 2015), thus leading to timing failures at low voltages. In other words, the clock skew degradation at low voltages typically defines V_{min} . Similar trends are observed when designing the clock network at low voltages and running at above-threshold voltages.

From the above considerations, the design of the clock network of integrated systems operating in a wide voltage range entails a fundamental tradeoff between the performance at above-threshold voltages, and the ability to scale down to low voltages. Various approaches have been proposed to make this tradeoff more favorable, and mitigate the skew-energy penalty imposed by the adoption of deep or shallow clock networks. For example, moderately deep networks with long-channel LVT buffers have been proposed in Myers et al. (2015). Design methodologies have been introduced in Seok et al. (2011), Tolbert et al. (2011), Zhao et al. (2012) to optimally design clock networks, although for a single low voltage. Techniques for adaptive point-to-point interconnects with regenerative drivers have been also proposed (Kim et al. 2014b; Wang et al. 2015), although they cannot be used for clock networks and are not supported by commercial EDA tools. Voltage-adaptive delay insertion across different clock domains was introduced in (Jain et al. 2012; Tokunaga et al. 2014) to mitigate the inter-domain skew (e.g., between processor and

memory), although no adaption to voltage has been performed within each clock domain. Clock network adaptation to a wide range of voltages with each clock domain has been demonstrated in Lin et al. (2017), where the clock network topology is reconfigured to minimize the skew at each specific voltage.

Regarding the supply network, voltage drops are less of a concern at near-threshold voltages and below, as the I_{on} transistor current is at least an order of magnitude lower than at nominal voltage. Accordingly, the current density drawn by the digital circuit is reduced by the same amount, and hence issues related to voltage drops across the supply network are largely mitigated. This partially alleviates the problem of the stronger impact of V_{DD} fluctuations at near-threshold voltages, due to the larger sensitivity of performance (see Sect. 4.7.2). This translates in a relaxed requirement on the supply rail width in the cell library, which can help slightly reduce the cell height. For analogous reasons, the lower clock frequency of near-threshold circuits makes the effect of the wire parasitic inductance negligible. Finally, the peak current absorbed by near-threshold circuits is also reduced by an order of magnitude, compared to above-threshold operation. Hence, the size of decoupling capacitors to keep V_{DD} fluctuations within a targeted band can be reduced by the same amount, thus saving area and improving the utilization factor of the module under design.

4.11 Perspectives and Trends

In summary, near-threshold circuits pose challenge and opportunities that are significantly differ from conventional above-threshold low-power circuits. Counteracting leakage in spite of the inefficacy of conventional leakage reduction techniques (see Sect. 4.5) requires a radically different approach that maximizes the opportunities to reduce leakage when transistors are not being used. This can be accomplished by introducing fine-grain power domains that can be power gated (e.g., with gate boosting to improve its effectiveness, as in Sect. 4.5) or voltage scaled to mitigate the leakage contribution of unused transistors. Power domains are typically coarse and of the size of at least an entire microprocessor, whereas such fine-grain power domains have the size of sub-blocks or execution units (e.g., ALU), or even finer (e.g., individual operators in the ALU). Although such approach certainly enhances the chances to turn off transistors, its direct application leads to significant area/energy/performance overhead. The latter is due to the need for additional power domain control circuitry, as well as isolation/clamping cells for power gating and level shifters (see Chap. 9) at the boundary of each domain.

Fine-grain voltage domains are also a highly promising approach in near-threshold circuits. Indeed, the ability to distribute different voltages with fine granularity maximizes the opportunities to correct variations in paths that turn out to be critical due to random variations, while reducing the energy in all other domains. The effectiveness of fine-grain voltage domains is further enhanced by the strong sensitivity of performance on V_{DD} (see Sect. 4.2), which ensures that voltage boosting is kept small (e.g., 100–200 mV) in all practical cases. For example, selective boosting can be used to reduce the general V_{min} of the circuit, while raising the voltage of the small portion of the circuit that needs to operate at higher voltages (Tokunaga et al. 2014). As another example, (Muramatsu et al. 2011) leverages such small voltage difference across voltage domains by suppressing level shifters

altogether, so that the voltages can be freely assigned to very small domains to compensate variations where they arise, while avoiding the otherwise large overhead of level shifters. The Panoptic approach (Putic et al. 2009) introduces both spatial and temporal fine granularity by using multiple sleep transistors that also dynamically connect sub-blocks to three different supply voltages. The sleep transistors serve the purpose of reducing leakage of unused sub-blocks, and assign them the minimum possible voltage for the task at hand when used.

Variations can also be exploited rather than added as design margin, when an adequately large number of replicas of a given block are available on the same chip. For example, Raghunathan et al. (2013) introduces the concept of “cherry picking” among many on-chip cores, which consists in the post-silicon selection of the most energy-efficient cores while keeping the others off. This permits to maximize the energy efficiency by leveraging the inevitable random variations, rather than tolerating them, at the cost of area due to the partial utilization of cores. Observe that full utilization would not be allowed anyway in practical cases, due to the “dark silicon” issue (see Chap. 1) that is determined by the chip power constraint.

In general, variations can be mitigated at different times, from design time to testing, chip boot time and run-time, as summarized in Fig. 4.54. At design time, all variation contributions need to be incorporated into the design (e.g., cycle time) margin, as they are not known upfront. The margin is lowered at testing time, as process variations are known and can hence be suppressed, whereas voltage, temperature and aging-induced variations are need to be included (as they will be defined later at in-field operation). At boot time, aging can be compensated as well. The margin is made very small and virtually removed when variations are compensated at run-time, i.e., when all process, temperature, (slow) voltage variations are known. Obviously, the cost of such detection and compensation of variations increases when moving from design to run-time.

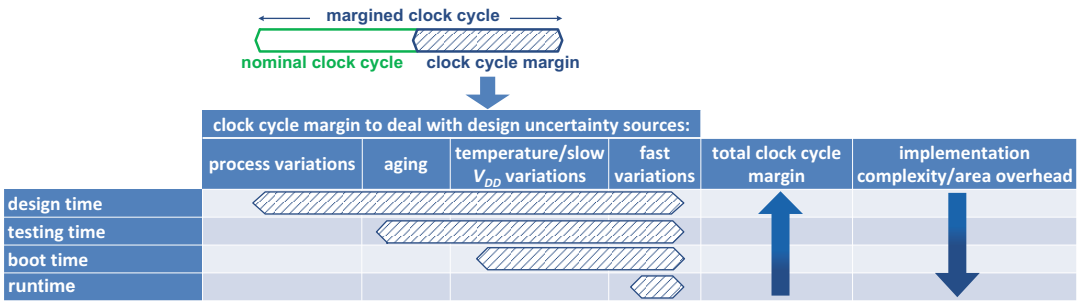


Fig. 4.54 Summary of techniques to counteract variations at different time, and resulting cycle margin and overhead

Due to very large design margin required at near-threshold voltages (see Sect. 4.7), adequate yield and energy efficiency certainly require the adoption of run-time compensation of variations. This is typically performed through timing error detection and correction (EDAC) methods, which have been investigated since early 2000s (Ernst et al. 2003), EDAC methods sense the timing margin at run time by detecting timing failures, so that the system can be tuned to operate at nearly-zero margin (Ernst et al. 2003). This permits to run at the highest possible frequency at given voltage, or at the minimum possible voltage at given frequency. Hence, error detection and correction improves the energy efficiency of circuits operating at any voltage, typically by 1.3–1.45X (see references below).

Error detection can be performed through canary circuits and Tunable Replica Circuits (TRCs) mimicking critical path variations, and hence predicting the occurrence of timing violations with high (but not 100%) level of confidence, at rather low overhead (Bowman et al. 2011). However, tracking the critical path across a wide range of voltages is difficult, and hence such methods are more appropriate for operation on a narrow range. Also, since TRCs try to replicate the critical path, they cannot completely eliminate the design margin. *In-situ* error detection is performed by inserting timing sensors to detect true timing failures, which typically entails significant area overhead. Several *in-situ* error detection methods have been proposed, such as Razor (Ernst et al. 2003), Razor II (Das et al. 2009), EDS (Bowman et al. 2009; Bowman et al. 2011), ERSA (Leem et al.

2010), Bubble Razor (Fojtik et al. 2013). However, their overhead is in the order of various (if not several) tens of percentage points, and hence an order of magnitude larger than TRCs, which has prevented their adoption in commercial chips. Recently, very low-overhead (i.e., percentage points) *in-situ* approaches have been demonstrated, such as Razor-Lite (Kwon et al. 2014) and iRazor (Zhang et al. 2016) for processors, and RazorSRAM for on-chip memories (Khayatzadeh et al. 2016). Being very lightweight, these approaches promise a much wider adoption of *in-situ* error detection in mass produced chips.

Another very promising direction to further reduce the energy per computation is offered by its tradeoff with quality. As discussed in Chap. 1, quality can be defined in different ways depending on the application and the sub-system under design. In a processing sub-system, quality is related to accuracy in terms of precision in case of arithmetic tasks, misclassification rate in classification tasks, or effective number of bits in an Analog-to-Digital Converter (ADC). The concept is far more general than approximations (e.g., approximate computing), in that it applies to a broad range of types of tasks and applications, and the tradeoff between quality and energy is dynamic and based on quality sensing (see example in Sect. 1.6.2).

Based on the concepts described in Sect. 1.6.2, energy-quality scalability has been introduced in many different sub-systems and levels of abstraction. For example, the first energy-quality SRAM memory has been

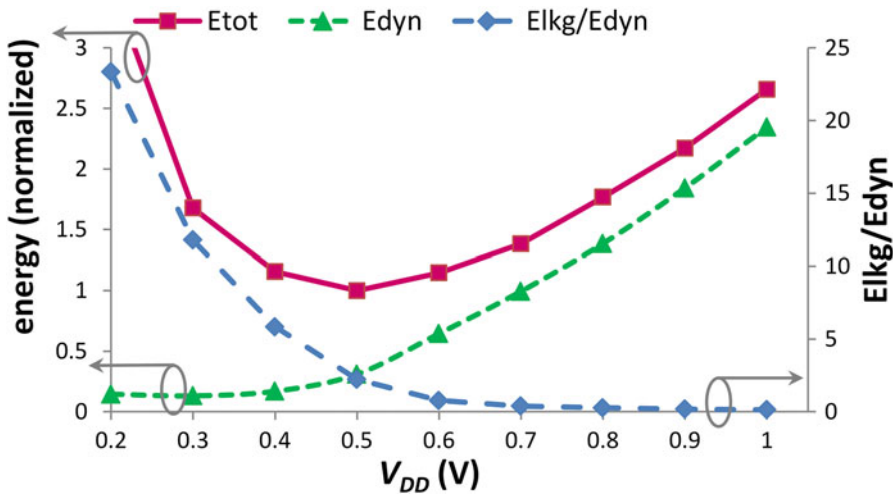


Fig. 4.55 Energy curve vs. V_{DD} in a 32-bit multiplier in 28-nm technology

introduced in Frustaci et al. (2015), where occasional faults (e.g., bitcells with inadequate write- or read-ability) occur in the array. The scalability comes from a bit-level management of the tradeoff between the bit error rate and the energy, by adjusting assist techniques (see Chap. 5) differently for different positions and in a dynamically scalable manner. This is beyond traditional memories where assist is uniformly applied to all bit positions and to fully suppress errors, which entails a substantial error cost. Similarly, selective Error Correction Codes have been introduced in the SRAM, to favor the robustness of the bits carrying the highest information content (e.g., MSBs in video processing applications), while saving on the other bit positions. Overall, this approach leads permits to improve the general quality by spending some energy in selected bit positions, thus enabling much more aggressive scaling on all positions and hence achieving quadratic benefit. Energy reductions of 2X have been demonstrated compared to traditional voltage scaling, at iso-quality (Frustaci et al. 2016). The same general concept has been applied to several other sub-systems, such as ADCs with dynamically scalable resolution (Freyman et al. 2014; Yip et al. 2011). In this case, when the application can tolerate a reduction in the ADC resolution, a more than 2X energy reduction is gained when

the resolution is reduced by one bit, leading to an exponential energy saving.

Finally, the presence of a minimum-energy point (MEP) actually poses a fundamental challenge in terms of energy scalability when the system is operating at the right of the MEP. Indeed, the MEP tends to be a flat minimum as discussed in Sect. 4.3.2, which in turns translates into insignificant energy savings when the voltage is scaled down from values at the right of the MEP towards the MEP itself. As an example, from Fig. 4.55 there is almost no energy saving when scaling from 0.6 V (i.e., at the right of the MEP) down to 0.5 V (i.e., the MEP), due to the flatness of the MEP. In other words, the voltage scalability of the design (i.e., its ability to operate at very low voltages) does not translate in an actual energy scalability (i.e., the ability to reduce energy when scaling down the voltage). To preserve energy scalability, the energy curve in Fig. 4.55 needs to be steep rather than flat, which is achieved only if the operating voltage is far enough on the right side of the MEP. For example, quadratic benefit is observed in this figure, at voltages from 0.8 V to 1 V. Conversely, to achieve good energy scalability at a given low voltage (e.g., 0.5 V), the MEP needs to be pushed to the left of this targeted voltage (e.g., 0.3 V). In other words, innovation is needed to move the MEP where needed, depending on the operating

voltage. At above-threshold voltages, the MEP can lie at a fairly high voltage, while not being a problem since the dominance of the dynamic energy still assured a quadratic benefit when downscaling V_{DD} . When a near-threshold voltage is targeted and further voltage scaling needs to be applied, the MEP needs to be dynamically moved to the left to make the energy curve steeper, and again achieve a nearly-quadratic energy benefit. We believe that this is one of the fundamental challenges that needs to be addressed to further improve the energy efficiency of low-voltage integrated systems for IoT.

Acknowledgement The authors acknowledge the kind support by the MOE2014-T2-2-158 grant from the Singaporean Ministry of Education.

References

- F. Abouzeid, S. Clerc, B. Pelloux-Prayer, F. Argoud, P. Roche, 28 nm CMOS, energy efficient and variability tolerant, 350 mV-to-1.0 V, 10 MHz/700MHz, 252bits frame error-decoder, in *Proceedings of ESSCIRC 2012* (Bordeaux, France, Sept. 2012), pp. 153–156
- M. Alioto, G. Scotti, A. Trifiletti, A novel framework to estimate the path delay variability via the fan-out-of-4 metric. *IEEE Trans. Circuits Syst.—Part I* (in press)
- M. Alioto, Understanding DC behavior of subthreshold CMOS logic through closed-form analysis. *IEEE Trans. Circuits Syst.—part I* **57**(7), 1597–1607 (2010)
- M. Alioto, Ultra-low power VLSI circuit design demystified and explained: a tutorial. *IEEE Trans. Circuits Syst.—part I* (invited) **59**(1), 3–29 (2012)
- M. Alioto, Challenges and techniques for ultra-low voltage logic with nearly-minimum energy. in *Short course at VLSI Symposium 2014*, Hawaii 10 June 2014
- M. Alioto, G. Palumbo, M. Pennisi, Understanding the effect of process variations on the delay of static and Domino logic. *IEEE Trans. VLSI Syst.* **18**(5), 697–710 (2010)
- M. Alioto, Guest editorial for the special issue on “Ultra-low-voltage VLSI circuits and systems for green computing. *IEEE Trans. Circuits Systems—part II* **59**(12), 849–852 (2012)
- M. Alioto, E. Consoli, G. Palumbo, *Flip-Flop Design in Nanometer CMOS—From High Speed to Low Energy* (Springer, Berlin, 2015)
- Z. Bo, D. Blaauw, D. Sylvester, K. Flautner, Theoretical and practical limits of dynamic voltage scaling, in *Proceedings of DAC* (2004), pp. 868–873
- K.A. Bowman, J.W. Tschanz, N.S. Kim, J.C. Lee, C.B. Wilkerson, S.-L. Lu, T. Karnik, V. De, Energy-efficient and metastability-immune resilient circuits for dynamic variation tolerance. *IEEE J. Solid-State Circuits* **44**, 49–63 (2009)
- K.A. Bowman, J.W. Tschanz, S.-L. Lu, P. Aseron, M. Khellah, A. Raychowdhury, B. Geuskens, C. Tokunaga, C. Wilkerson, T. Karnik, V. De, A 45 nm resilient microprocessor core for dynamic variation tolerance. *IEEE J. Solid-State Circuits* **46**(1), 194–208 (2011)
- T. Burd, T. Pering, A. Stratakos, R. Brodersen, A dynamic voltage scaled microprocessor system, in *IEEE ISSCC Dig. Tech. Papers* (Feb. 2015), pp. 294–295
- A. Chandrakasan, D. Daly, D. Finchelstein, J. Kwong, Y. Ramadass, M. Sinangil, V. Sze, N. Verma, Technologies for ultradynamic voltage scaling. *Proc. IEEE* **98**(2), 191–214 (2010)
- D. Chinnery, K. Keutzer, *Closing the Power Gap between ASIC & Custom* (Springer, Berlin, 2007)
- J. Crop, E. Krimer, N. Moezzi-Madani, R. Pawlowski, T. Ruggeri, P. Chiang, M. Erez, Error detection and recovery techniques for variation-aware CMOS computing: a comprehensive review. *J. Low Power Electron. Appl.* **1**, 334–356 (2011)
- S. Das, C. Tokunaga, S. Pant, W.-H. Ma, S. Kalaiselvan, K. Lai, D.M. Bull, D.T. Blaauw, Razor II: in situ error detection and correction for PVT and SER tolerance. *IEEE J. Solid-State Circuits* **44**(1), 32–48 (2009)
- G. De Micheli, *Synthesis and Optimization of Digital Circuits* (McGraw Hill, New York, 1994)
- R. Dreslinski, M. Wiecekowski, D. Blaauw, D. Sylvester, T. Mudge, Near-threshold computing: reclaiming Moore’s law through energy efficient integrated circuits. *Proc. IEEE* **98**(2), 253–266 (2010)
- C. Enz, E. Vittoz, *Charge-Based MOS Transistor Modeling: The EKV Model for Low-Power and RF IC Design* (Wiley, New York, 2006)
- D. Ernst, N.S. Kim, S. Das, S. Pant, R. Rao, T. Pham, C. Ziesler, D. Blaauw, T. Austin, K. Flautner, T. Mudge, Razor: a low-power pipeline based on circuit-level timing speculation, in *Proceedings of MICRO-36* (Dec. 2003), pp. 7–18
- D. Flynn, R. Aitken, A. Gibbons, K. Shi, *Low Power Methodology Manual* (Springer, New York, 2007)
- M. Fojtik, D. Fick, Y. Kim, N. Pinckney, D. Harris, D. Blaauw, D. Sylvester, Bubble razor: eliminating timing margins in an ARM Cortex-M3 processor in 45 nm CMOS using architecturally independent error detection and correction. *IEEE J. Solid-State Circuits* **48**(1), 66–81 (2013)
- L. Freyman, D. Fick, M. Alioto, D. Blaauw, D. Sylvester, A 346 μm^2 VCO-based, reference-free, self-timed sensor interface for cubic-millimeter sensor nodes in 28 nm CMOS. *IEEE J. Solid-State Circuits* **49**(11), 2462–2473 (2014)
- F. Frustaci, M. Khayatzaeh, D. Blaauw, D. Sylvester, M. Alioto, SRAM for error-tolerant applications with dynamic energy-quality management in 28 nm CMOS. *IEEE J. Solid-State Circuits* **50**(3), 1310–1323 (2015)

- F. Frustaci, D. Blaauw, D. Sylvester, M. Alioto, Approximate SRAMs with dynamic energy-quality management. *IEEE Trans. VLSI Syst.* **24**(6), 2128–2141 (2016)
- G. Gammie, N. Ickes, M. Sinangil, R. Rithe, J. Gu, A. Wang, H. Mair, S. Datla, R. Bing, S. Honnavara-Prasad, L. Ho, G. Baldwin, D. Buss, A. Chandrakasan, U. Ko, A 28 nm 0.6 V low-power DSP for mobile applications, in *ISSCC Digest of Technical Papers (ISSCC)* (San Francisco, Feb. 2011)
- T. Gemmeke, M. Ashouei, B. Liu, M. Meixner, T.G. Noll, H. de Groot, Cell libraries for robust low-voltage operation in nanometer technologies. *Solid-State Electron.* **84**, 132–141 (2013)
- J. Gregg, T.W. Chen, Post silicon power/performance optimization in the presence of process variations using individual well-adaptive body biasing. *IEEE Trans. VLSI Syst.* **15**(3), 366–376 (2007)
- S. Hanson, B. Zhai, D. Blaauw, D. Sylvester, A. Bryant, X. Wang, Energy optimality and variability in sub-threshold design, in *Proceedings of ISLPED 2006* (2006), pp. 363–365
- S. Hanson, B. Zhai, K. Bernstein, D. Blaauw, A. Bryant, L. Chang, K.K. Das, W. Haensch, E.J. Nowak, D.M. Sylvester, Ultralow-voltage, minimum-energy CMOS. *IBM J. Res. & Dev.* **50**(4/5) (2006), pp. 469–490
- S. Hanson, B. Zhai, M. Seok, B. Cline, K. Zhou, M. Singhal, M. Minuth, J. Olson, L. Nazhandali, T. Austin, D. Sylvester, D. Blaauw, Exploring variability and performance in a sub-200-mV processor. *IEEE J. Solid-State Circuits* **43**(4), 881–890 (2008)
- D. Harris, R. Ho, G.-Y. Wei, M. Horowitz, The Fanout-of-4 inverter delay metric, unpublished manuscript <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.68.831&rep=rep1&type=pdf>
- S. Hsu, A. Agarwal, M. Anders, S. Mathew, H. Kaul, F. Sheikh, R. Krishnamurthy, A 280 mV-to-1.1 V 256b reconfigurable SIMD vector permutation engine with 2-dimensional shuffle in 22 nm CMOS, in *ISSCC Digest of Technical Papers (ISSCC)* (San Francisco, Feb. 2012)
- International Technology Roadmap for Semiconductors: 2013 edition. <http://www.itrs.net> (2013)
- D. Jacquet et al., 2.6GHz ultra-wide voltage range energy efficient dual A9 in 28 nm UTBB FD-SOI, in *IEEE Symposium on VLSI Circuits Dig. Tech. Papers* (June 2013)
- S. Jain et al., A 280 mV-to-1.2 V wide-operating-range IA-32 processor in 32 nm CMOS, in *IEEE ISSCC Dig. Tech. Papers* (Feb. 2012), pp. 66–67
- D. Jeon, M. Seok, C. Chakrabarti, D. Blaauw, D. Sylvester, A super-pipelined energy efficient sub-threshold 240 MS/s FFT core in 65 nm CMOS. *IEEE J. Solid-State Circuits* **47**(1), 23–34 (2013)
- H. Kaul, M.A. Anders, S.K. Mathew, S.K. Hsu, A. Agarwal, F. Sheikh, R.K. Krishnamurthy, S. Borkar, A 1.45 GHz 52-to-162GFLOPS/W variable-precision floating-point fused multiply-add unit with certainty tracking in 32 nm CMOS, in *IEEE ISSCC Dig. Tech. Papers*, (Feb. 2012), pp. 182–183
- M. Khayatzadeh, M. Saligane, J. Wang, M. Alioto, D. Blaauw, D. Sylvester, A reconfigurable dual port memory with error detection and correction in 28 nm FDSOI, in *IEEE ISSCC Dig. Tech. Papers* (Feb. 2016), pp. 310–311
- Y. Kim, W. Jung, I. Lee, Q. Dong, M. Henry, D. Sylvester, D. Blaauw, A static contention-free single-phase-clocked 24T Flip-Flop in 45 nm for low-power applications. in *IEEE ISSCC Dig. Tech. Papers* (Feb. 2014)
- S. Kim, M. Seok, Reconfigurable interconnect-driving technique for ultra-dynamic-voltage-scaling systems, in *IEEE ACM International Symposium on Low Power Electronics and Design (ISLPED)* (2014)
- I. Kwon, S. Kim, D. Fick, M. Kim, Y.-P. Chen, D. Sylvester, Razor-lite: a light-weight register for error detection by observing virtual supply rails. *IEEE J. Solid-State Circuits* **49**(9), 2054–2066 (2014)
- L. Leem, H. Cho, J. Bau, Q.A. Jacobson, S. Mitra, ERSA: error resilient system architecture for probabilistic applications, in *Proceedings of DATE 2010* (Dresden, Germany, Mar. 2010), pp. 1560–1565
- L. Lin, S. Jain, M. Alioto, Reconfigurable clock networks for random skew mitigation from sub-threshold to nominal voltage, in *IEEE ISSCC Dig. Tech. Papers* (Feb. 2017)
- M. Alioto, G. Scotti, A. Trifiletti, A novel framework to estimate the path delay variability via the Fan-Out-of-4 metric. *IEEE Trans. Circuits Syst.—part I*
- D. Markovic, V. Stojanovic, B. Nikolic, M.A. Horowitz, R.W. Brodersen, Methods for true energy-performance optimization. *IEEE J. Solid-State Circuits* **39**(8), 1282–1293 (2004)
- S.M. Martin, K. Flautner, T. Mudge, D. Blaauw, Combined dynamic voltage scaling and adaptive body biasing for lower power microprocessors under dynamic workloads, in *Proceedings of ICCAD'02* (Nov. 2002), pp. 721–725
- M. Meijer, J. Pineda de Gyvez, Body-bias-driven design strategy for area- and performance-efficient CMOS circuits. *IEEE Trans. VLSI Syst.* **20**(1), 42–51 (2012)
- M. Merrett, Y. Wang, M. Alioto, M. Zwolinski, Design metrics for RTL level estimation of delay variability due to intradie (random) variations, in *Proceedings of ISCAS 2010* (Paris (France), May 2010), pp. 2498–2501
- A. Muramatsu, T. Yasufuku, M. Nomura, M. Takamiya, H. Shinohara, T. Sakurai, 12% power reduction by within-functional-block fine-grained adaptive dual supply voltage control in logic circuits with 42 voltage domains, in *37th European Solid-State Circuits Conference (ESSCIRC)* (Helsinki (Finland), Sep. 2011), pp. 191–194
- J. Myers, A. Savanth, D. Howard, R. Gaddh, P. Prabhat, D. Flynn, An 80nW retention 11.7pJ/cycle active sub-threshold ARM Cortex®-M0+ sub-system in

- 65 nm CMOS for WSN applications, in *IEEE ISSCC Dig. Tech. Papers* (Feb. 2015), pp. 144–145
- S. Narayanan, J. Sartori, R. Kumar, D.L. Jones, Scalable stochastic processors, in *Proceedings of DATE 2010* (Dresden, Germany, Mar. 2010), pp. 335–338
- S. Narendra, A. Chandrakasan (eds.), *Leakage in Nanometer CMOS Technologies* (Springer, Berlin, 2006)
- K. Nose, T. Sakurai, Optimization of VDD and VTH for low-power and high-speed applications, in *Proceedings of ASPDAC* (Jan. 2000), pp. 469–474
- M. Olivieri, G. Scotti, A. Trifiletti, A novel yield optimization technique for digital CMOS circuits design by means of process parameters run-time estimation and body bias active control. *IEEE Trans. VLSI Syst.* **13** (5), 630–638 (2005)
- M.M. Orshansky, S. Nassif, D. Boning, *Design for Manufacturability and Statistical Design* (Springer, Berlin, 2008)
- D. Patil, M. Horowitz, Joint supply, threshold voltage and sizing optimization for design of robust digital circuits. <http://vlsiweb.stanford.edu/papers/JointVddVthSizing.pdf>
- M. Pelgrom, A. Duijnmaier, A. Welbers, Matching properties of MOS transistors. *IEEE J. Solid-State Circuits* **24**(1), 1433–1439 (1989)
- M. Putic, L. Di, B. H. Calhoun, J. Lach, Panoptic DVS: a fine-grained dynamic voltage scaling framework for energy scalable CMOS design, in *Proceedings of ICCD 2009* (Lake Tahoe, CA, Oct. 2009), pp. 491–497
- B. Raghunathan, Y. Turakhia, S. Garg, D. Marculescu, Cherry-picking: exploiting process variations in dark-silicon homogeneous chip multi-processors, in *Proceedings of DATE 2013* (Grenoble, France, Mar. 2013), pp. 39–44
- Y.K. Ramadass, A.P. Chandrakasan, Minimum energy tracking loop with embedded DC–DC converter enabling ultra-low-voltage operation down to 250 mV in 65 nm CMOS. *IEEE J. Solid-state Circuits* **43**(1), 256–265 (2008)
- T. Sakurai, R. Newton, Alpha-power law MOSFET model and its applications to CMOS inverter delay and other formulas. *IEEE J. Solid-State Circuits* **25**(2), 584–594 (1990)
- W. Sansen, *Analog Design Essentials* (Springer, New York, 2006)
- M. Seok, D. Blaauw, D. Sylvester, Robust clock network design methodology for ultra-low voltage operations, in *IEEE Transactions on Emerging Selected Topics Circuits Systems*, vol. 1(2) (2011)
- F. Sheikh, S. Mathew, M. Anders, H. Kaul, S. Hsu, A. Agarwal, R. Krishnamurthy, S. Borkar, A 2.05 GVertices/s 151 mW lighting accelerator for 3D graphics vertex and pixel shading in 32 nm CMOS, in *IEEE ISSCC Dig. Tech. Papers* (Feb. 2012), pp. 178–179
- V. Srinivasan, D. Brooks, M. Gschwind, P. Bose, V. Zyuban, P.N. Strenski, P.G. Emma, Optimizing pipelines for power and performance, in *Proceedings of International Symposium on Microarchitectures* (2002), pp. 333–344
- I. Sutherland, B. Sproull, D. Harris, *Logical Effort: Designing Fast CMOS Circuits* (Morgan-Kaufmann, Burlington, 1999)
- C. Tokunaga, J.F. Ryan, C. Augustine, J.P. Kulkarni, Y.-C. Shih, S.T. Kim, R. Jain, K. Bowman, A. Raychowdhury, M.M. Khellah, J.W. Tschanz, V. De, A graphics execution core in 22 nm CMOS featuring adaptive clocking, selective boosting and state-retentive sleep, in *ISSCC Digest of Technical Papers (ISSCC)* (San Francisco, CA, Feb. 2014)
- J.R. Tolbert, X. Zhao, S.K. Lim, S. Mukhopadhyay, Analysis and design of energy and slew aware subthreshold clock systems. *IEEE Trans. CAD* **30**(9), 1348–1358 (2011)
- J.W. Tschanz, J.T. Kao, S.G. Narendra, R. Nair, D.A. Antoniadis, A.P. Chandrakasan, V. De, Adaptive body bias for reducing impacts of die-to-die and within-die parameter variations on microprocessor frequency and leakage. *IEEE J. Solid-State Circuits* **37**(11), 1396–1402 (2002)
- Y. Tsidividis, *Operational Modeling of the MOS Transistor*, 2nd edn. (McGraw-Hill, New York, 1999)
- R.E. Walpole, R.H. Myers, S.L. Myers, K. Ye, *Probability & Statistics for Engineers & Scientists* (Prentice Hall, Englewood Cliffs, 2006)
- A. Wang, A. Chandrakasan, 180-mV subthreshold FFT processor using a minimum energy design methodology. *IEEE J. Solid-State Circuits* **40**(1), 310–319 (2005)
- A. Wang, B.H. Calhoun, A. Chandrakasan, *Sub-threshold design for ultra low-power systems* (Springer, Berlin, 2006)
- J. Wang, N. Pinckney, D. Blaauw, D. Sylvester, Reconfigurable self-timed regenerators for wide-range voltage scaled interconnect, in *Proceedings of ASSCC 2015* (Nov. 2015)
- N. Weste, D. Harris, *CMOS VLSI Design*, 4th edn. (Pearson Education, Upper Saddle River, 2011)
- R. Wilson et al., A 460 MHz at 397 mV, 2.6 GHz at 1.3 V, 32b VLIW DSP, embedding FMAX tracking, in *IEEE ISSCC Dig. Tech. Papers* (Feb. 2014), pp. 452–453
- T. Xanthopoulos, *Clocking in Modern VLSI Systems* (Springer, New York, 2009)
- M. Yip, A. Chandrakasan, A resolution-reconfigurable 5-to-10 b 0.4-to-1 V power scalable SAR ADC, in *IEEE ISSCC Dig. Tech. Papers* (2011), pp. 190–191
- Y. Zhang, M. Khayatzaeh, K. Yang, M. Saligane, M. Alioto, D. Blaauw, D. Sylvester, iRazor: 3-transistor current-based error detection and correction in an ARM Cortex-R4 processor, in *IEEE ISSCC Dig. Tech. Papers* (Feb. 2016), pp. 160–161
- X. Zhao, J.R. Tolbert, S. Mukhopadhyay, S.K. Lim, Variation-aware clock network design methodology for ultralow voltage (ULV) circuits. *IEEE Trans. CAD* **31**(8), 1222–1234 (2012)
- W. Zhao, Y. Ha, M. Alioto, Novel self-body-biasing and statistical design for near-threshold circuits with ultra energy-efficient AES as case study. *IEEE Trans. VLSI Syst.* **23**(8), 1390–1401 (2015)