

Chapter 2

Analog-to-Digital Conversion

Fundamentals



The tremendous popularity but also challenges of data converters as key interface functions between the physical (analog) world and the electronic (digital) world were discussed from a bird's eye view in the previous chapter. Before delving into advanced architectural and design details, this chapter will cover the fundamental A/D conversion principles, some important performance metrics, as well as practical limitations, serving as the foundation for the following chapters.

Section 2.1 serves as a theoretical background by reviewing the two main functions in every A/D conversion: (1) sampling and (2) quantization. The major error sources stemming from the individual blocks of practical converters are identified and analyzed in Sect. 2.2, followed by a review of the most important performance metrics and figures of merit in Sect. 2.3. Section 2.4 derives the impact on the accuracy-speed-power for every major error source. This derivation leads to the establishment of the fundamental limits on a converter's performance, imposed by circuits, by technology, and ultimately by physics. The limits in this chapter form the basis of what may be theoretically achievable and, together with the architectural overheads presented in Chap. 3, serve as guidelines to assist the design choices of the prototypes in Chaps. 4–7. This chapter closes with an overview and conclusions in Sect. 2.5.

2.1 Theoretical Background

As already mentioned, every analog signal is continuous both in time and in amplitude. Therefore, two main processes are essential to obtain the final digital waveform:

1. Sampling (to achieve the time discretization)
2. Quantization (to achieve the amplitude discretization)

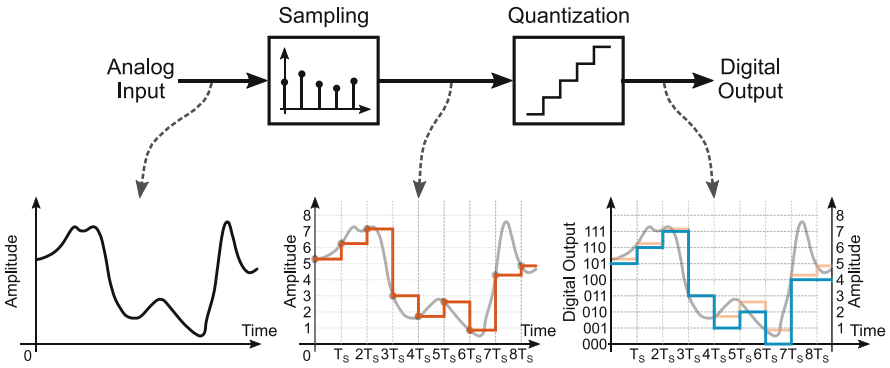


Fig. 2.1 Block diagram of an ideal A/D conversion (top) and the resulting waveforms at every part of the chain (bottom)

Figure 2.1 depicts the block diagram of an ideal Analog-to-Digital (A/D) conversion with its corresponding waveforms. The continuous time and amplitude analog input signal (black waveform) is uniformly sampled with a period of T_s (or at a sample rate¹ of f_s). The resulting time-discrete analog signal (orange waveform) updates its value only at integer multiples of T_s . When the time is equal to an integer multiple of T_s , the sampled signal is equal in value to the analog input at that instant and keeps its value until the next multiple of T_s arrives. Between two consecutive time instants, the sampled signal is held constant and can be further processed down the conversion chain.

Next, the quantization takes place, where the sampled signal is discretized in amplitude and its analog values are mapped onto a set of discrete levels (blue waveform). The digital output, now discrete in both time and amplitude, is an approximation of the initial analog input, with its approximation accuracy limited by the number of the available discrete levels. During both sampling and quantization, there is information loss since an error is introduced on the initial analog signal. This error can be reduced by increasing the number of time samples and/or the number of discrete levels. As we will see in the remainder of this book, guaranteeing simultaneously both can be far from trivial.

2.1.1 Sampling

Sampling is the basic process that transfers a waveform from the continuous time to the discrete time domain. The sampling process can be described mathematically

¹ Throughout this book, the terms sample rate, sampling rate, sampling speed, and/or sampling frequency will all refer to the same quantity f_s .

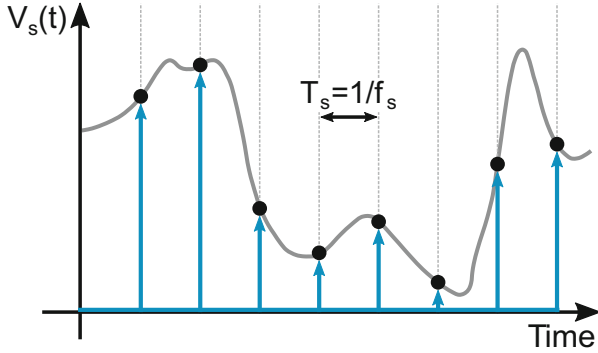


Fig. 2.2 Sampling a continuous-time signal using a Dirac pulse sequence

by means of the Dirac function $\delta(t)$, whose integral is equal to one at the integration instant and zero elsewhere [12]. The required sampling time frame is determined by a sequence of equidistant in time Dirac pulses, spaced by T_s . The time-discrete signal is a result of the multiplication of the Dirac pulses with the original waveform, with an amplitude equal to the amplitude of the waveform in the sampling instants and undefined elsewhere (Fig. 2.2). The mathematical formula expressing the above is given as

$$V_s(t) = V(t) \cdot \sum_{n=-\infty}^{n=\infty} \delta(t - nT_s) = \sum_{n=-\infty}^{n=\infty} V(nT_s). \quad (2.1)$$

Generally, the transformation of a signal from time domain to frequency domain is done by means of its Fourier Transform (FT). For a time-discrete signal specifically, this transformation in the frequency domain occurs by employing the signal's Discrete Fourier Transform (DFT). Taking into account that a multiplication in time is a convolution in frequency, the spectrum of the time-discrete signal $V_s(t)$ is depicted in Fig. 2.3 and given by

$$V_s(f) = \frac{1}{T_s} \cdot \sum_{n=-\infty}^{n=\infty} V(f - nf_s). \quad (2.2)$$

The dual-sided band around zero with a frequency content within $\pm f_{in}$ is attributed to the original waveform. The replica or alias bands around multiples of f_s result from the multiplication of the original waveform with the repetitive by $T_s = 1/f_s$ Dirac pulse sequence. The signal bands with the same frequency content around any multiple of f_s , after processing the spectrum with a Fast Fourier Transform (FFT) algorithm become indistinguishable around zero. As a numerical example, single-tone signals with 211 MHz, 789 MHz, 1.211 GHz, 1.789 GHz, and 2.211 GHz input

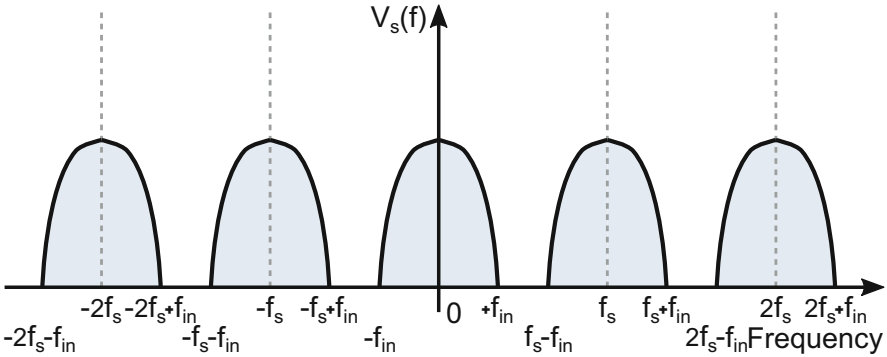


Fig. 2.3 Frequency spectrum of a signal multiplied with a sequence of Dirac pulses

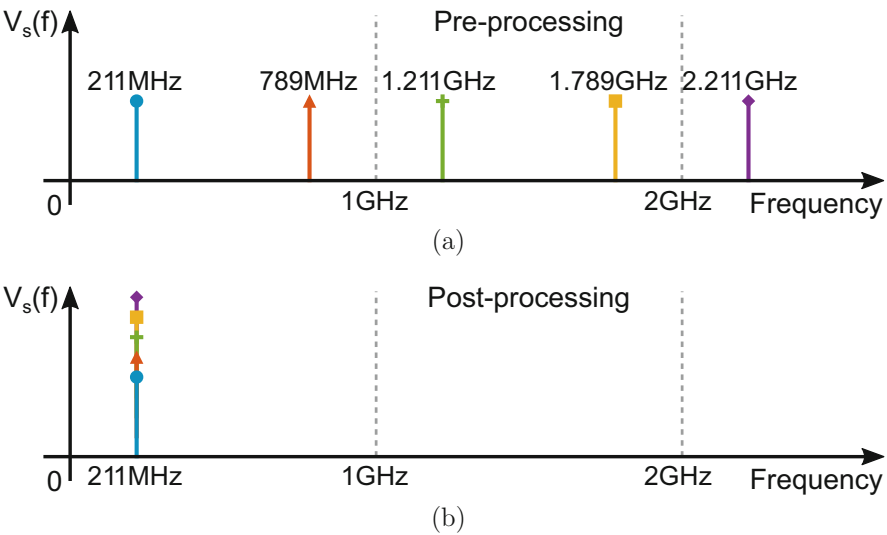


Fig. 2.4 (a) Single-tone signals with different frequencies (b) fall in the same frequency location after spectrum processing

frequencies (Fig. 2.4a) will all end up at the 211 MHz frequency location when sampled at 1 GS/s (Fig. 2.4b).

If the band of the original waveform increases in width, so will its alias bands. This will eventually lead to the bands overlapping, causing mixing of information between them and making it impossible to isolate the information from each band correctly. This irreversible situation is described as aliasing. In order to prevent information loss due to aliasing and yield the sampling process reversible, the following condition between the instantaneous signal bandwidth $f_{in,bw}$ and the sample rate f_s must be obeyed:

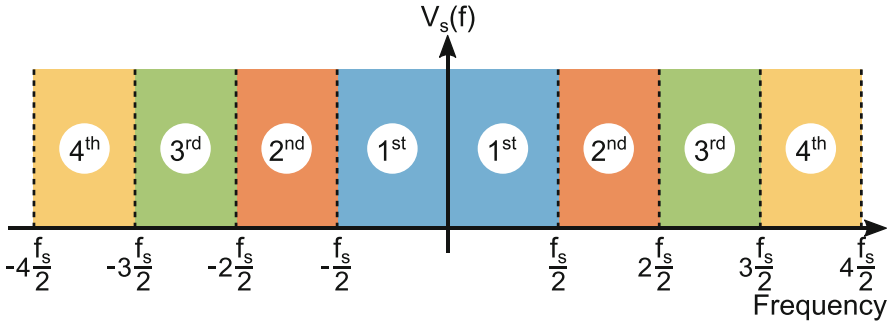


Fig. 2.5 Dual-sided frequency spectrum highlighting different Nyquist zones

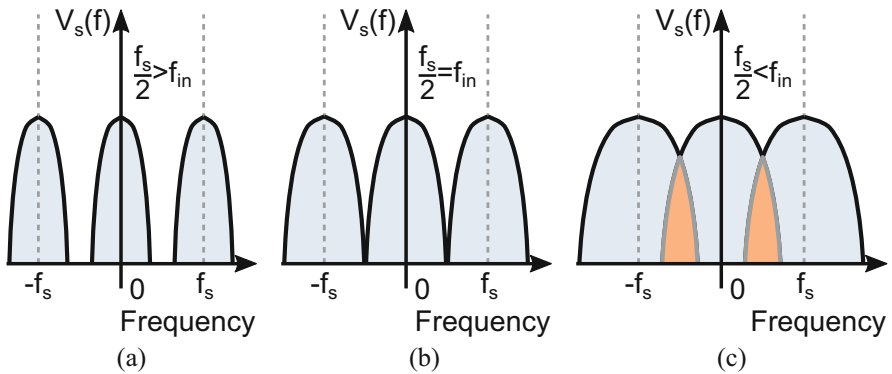


Fig. 2.6 (a), (b) Two cases of signals with bands meeting the Nyquist criterion and (c) one scenario where bands are overlapping leading to information loss

$$\frac{f_s}{2} > f_{in,bw}. \tag{2.3}$$

Known as the sampling theorem or Nyquist sampling criterion [13, 14], the above expression can be translated to

A band-limited continuous-time signal can be sampled and perfectly reconstructed if the sample rate is more than twice the signal’s instantaneous bandwidth. The frequency band between zero and $f_s/2$ is defined as the Nyquist bandwidth or the 1st Nyquist zone. The total spectrum comprises an infinite number of Nyquist zones, each with a width of $f_s/2$. Figure 2.5 shows the first four Nyquist zones in the spectrum, indicating their frequency allocation and width. For signals originally residing in the odd-order zones, their bands after sampling are copied to the 1st Nyquist zone as they are, while the bands of even-order zones are mirrored. Under the assumption that Eq. (2.3) holds (Fig. 2.6a, b), the original signal can be accurately reconstructed by a reconstruction filter. However, a violation

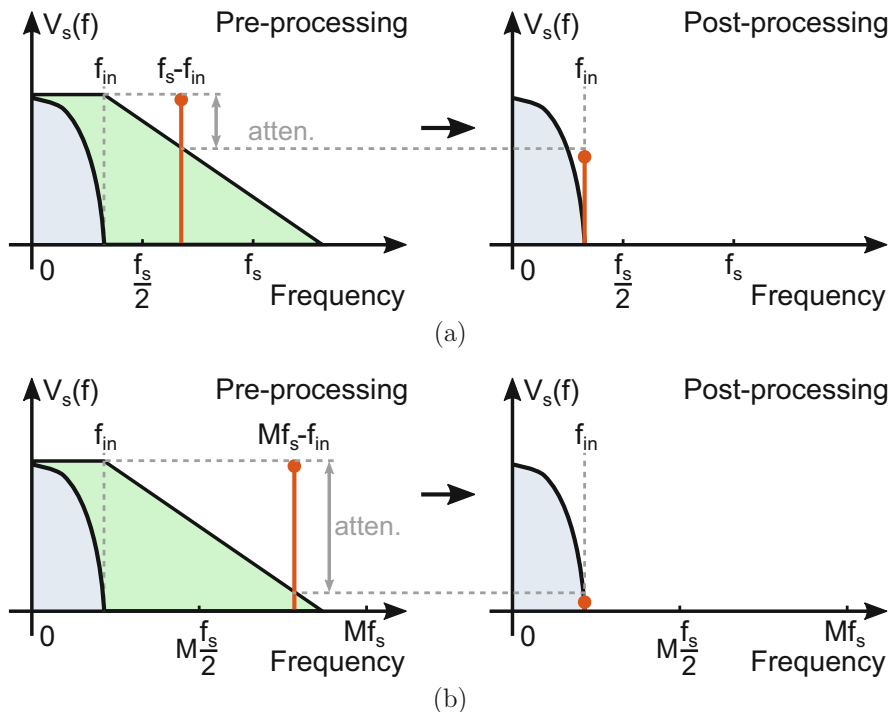


Fig. 2.7 Anti-aliasing filter on a parasitic tone when (a) sampling at Nyquist rate (slightly oversampled in practice) and (b) oversampling by $M > 1$

of the Nyquist criterion (Fig. 2.6c) will result in aliasing and render an accurate reconstruction of the original signal impossible.

Even if the useful signal resides within the Nyquist bandwidth, different types of undesired signals or interferers may appear at higher Nyquist zones, mixing up with the useful signal after sampling in the 1st Nyquist zone. Examples of such undesired signals are harmonic-related products of the main signal and/or interferers/noise from parts in the signal chain preceding the sampling. To prevent these unwanted signals from limiting the Dynamic Range (DR) of the chain, an anti-aliasing filter is typically employed prior to sampling to remove any component outside the Nyquist bandwidth. The specifications of this filter, whose implementation may include active and/or passive components, heavily depend on how much attenuation it needs to provide at which frequency distance with respect to $f_s/2$. Given that typical filters provide an attenuation of 20 dB/decade per order, a multi-order robust filter design becomes increasingly challenging and expensive as the frequency band of interest approaches $f_s/2$. Figure 2.7a illustrates the case of attenuating a parasitic tone by a finite-order anti-aliasing filter for a signal with $f_{in} < f_{in,bw}$ sampled at the Nyquist rate.

One way to improve the filter attenuation for a certain order or relax the filter order for a certain attenuation is to sample faster than the Nyquist criterion imposes. Increasing the sample rate (oversampling) provides a trade-off between parasitic tone attenuation and clock speed to sample and process data for a certain filter order [15]. Figure 2.7b illustrates how oversampling by a factor of M significantly improves the parasitic tone attenuation for the same filter order. However, for very wideband signals, generating the clock for a certain oversampling becomes equivalently challenging as increasing the filter order.

As a final note on sampling, it is worth mentioning that the Nyquist criterion is still satisfied and aliasing is not an issue for a signal residing in any of the Nyquist zones, as long as it is band-limited within one. In fact, this sampling property is utilized in the increasingly popular sub-sampling ADCs in communication systems. Directly sampling Intermediate Frequency (IF)/Radio Frequency (RF) signals in higher Nyquist zones and processing them digitally allow simplification of the signal chain by eliminating several frequency down-conversion blocks, such as a mixer, an IF amplifier, and filters. However, this increases the sub-sampling Analog-to-Digital Converter (ADC)'s bandwidth and spectral purity requirements at higher Nyquist zones. Chapter 6 of this book introduces circuit and architecture techniques for efficiently realizing wideband RF sampling ADCs.

2.1.2 Ideal Quantization

An ideal quantizer is a memoryless non-linear block, which uses B bits to translate the sampled signal to a digital word of binary format (0s and 1s). B represents the aggregate resolution with which the digital output resembles the analog input. Figure 2.8 shows the conceptual model and transfer characteristic of an ideal B -bit quantizer. Each signal value is compared against 2^B discrete levels, and its amplitude is rounded to the nearest level. The output Encoding Logic (ENC) decides how the rounding is done. The maximum input amplitude is defined as the Full-Scale (FS), and the difference between two adjacent transition levels (a.k.a. the step width), Δ , is quantified in the analog domain as the Least Significant Bit (LSB) such that $\Delta = \text{FS}/2^B$.

The digital word can be back-converted to a discrete amplitude analog signal V_q by multiplying each bit with its assigned binary weight, provided that the analog value of Δ is known

$$V_q = \Delta * \left(\sum_{i=0}^{B-1} bit_i * 2^i + bit_1 * 2^1 + bit_2 * 2^2 + \dots + bit_{B-1} * 2^{B-1} \right). \quad (2.4)$$

Due to the rounding process, there is a quantization error ϵ_q added to the original signal V_{in} , with a value ideally within $\pm\Delta/2$ for signals inside FS, while growing out of bounds outside FS (Fig. 2.8). The minimum error power is achieved for uniformly

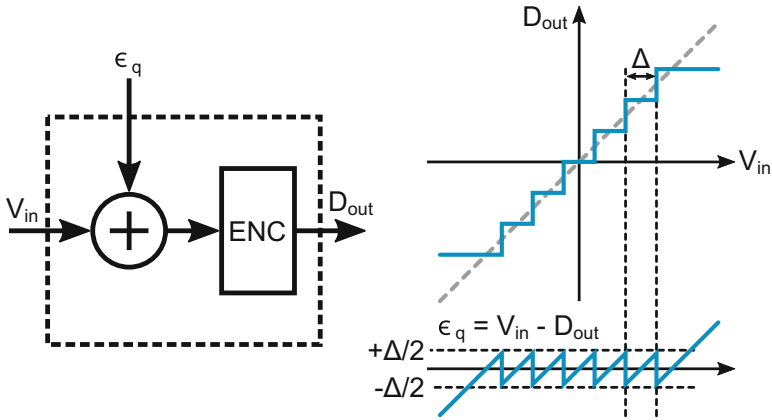
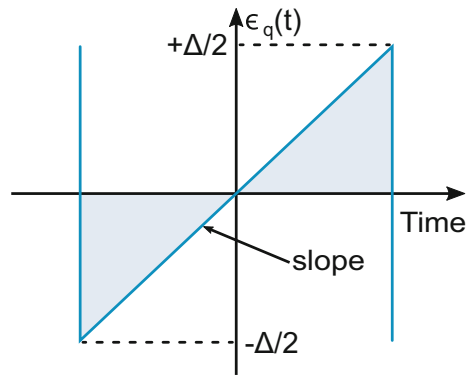


Fig. 2.8 Conceptual model and transfer characteristic of an ideal quantizer

Fig. 2.9 Sawtooth approximation of ϵ_q as a function of time



spaced discrete levels [16]. The back-converted signal relation with the original signal is expressed as

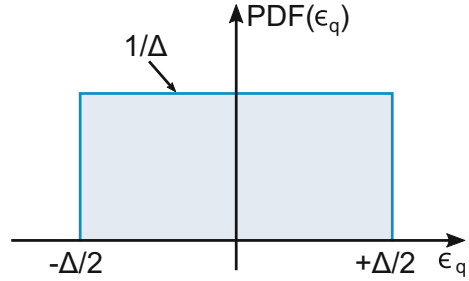
$$V_q = V_{in} + \bar{\epsilon}_q. \tag{2.5}$$

Strictly speaking, ϵ_q is a deterministic quantity, heavily depending on the properties of the signal at hand. For a linear ramp signal that contains several LSBs, ϵ_q can be approximated in time domain by a sawtooth waveform with a peak-to-peak amplitude of Δ , as shown in Fig. 2.9

$$\epsilon_q(t) = slope \cdot t, \quad -\frac{\Delta}{2} \leq slope \cdot t \leq \frac{\Delta}{2}. \tag{2.6}$$

Due to the signal periodicity, an integration over a single period of T_p suffices to calculate the Root-Mean-Square (RMS) value of the error

Fig. 2.10 Uniformly distributed PDF of ϵ_q within $\pm\Delta/2$



$$\bar{\epsilon}_q^2 = \frac{1}{T_p} \int_{-\frac{T_p}{2}}^{\frac{T_p}{2}} \epsilon_q^2(t) dt = \frac{\text{slope}}{\Delta} \int_{-\frac{T_p}{2}}^{\frac{T_p}{2}} \left(\frac{\Delta}{T_p}\right)^2 t dt = \frac{\Delta^2}{12} \Rightarrow \bar{\epsilon}_q = \frac{\Delta}{\sqrt{12}}. \quad (2.7)$$

If a more statistical approach is followed, considering that over a long time span all values within $\pm\Delta/2$ will show up with the same probability, ϵ_q assumes a uniform Probability Density Function (PDF) within that same region as is illustrated in Fig. 2.10. The necessary conditions for the validity of this approach are:

- The signal is sufficiently large or the quantizer resolution is large, such as to cover an adequate amount of levels
- The input is uncorrelated with the quantization error or the input frequency is not harmonically linked to the sample rate
- The signal is limited to FS, such that there is no quantizer overloading

If the above conditions hold, ϵ_q may be allocated a zero mean μ_{ϵ_q} and a variance $\sigma_{\epsilon_q}^2$ that can be calculated as in [17]

$$\sigma_{\epsilon_q}^2 = \bar{\epsilon}_q^2 = \frac{1}{\Delta} \int_{-\frac{\Delta}{2}}^{\frac{\Delta}{2}} \epsilon_q^2 d\epsilon_q = \frac{\Delta^2}{12}, \quad (2.8)$$

which matches the result of Eq. (2.7). As pointed out in [17], this quantization “noise” upon sampling shows a uniform spread across the entire Nyquist bandwidth. In case the input frequency is harmonically linked to the sample rate, there exists a relation between the input and ϵ_q resulting in the energy being accumulated in the harmonics of the signal. When performing a spectral analysis through FFT, this correlation can be avoided by choosing an integer number of signal periods (coherent sampling) and relatively prime number of periods and points [18]. Appendix A describes such an FFT setup.

As the quantizer resolution decreases, the non-linear nature of the quantization process dominates over its noise-like approximation, resulting in a distortion dominated spectrum rather than the flat noise-like. Figure 2.11 plots the spectra of an ideally quantized 77 MHz input signal coherently sampled at 1 GS/s for various resolutions. A reduction of about 8–9 dB per added bit is seen in the odd harmonic

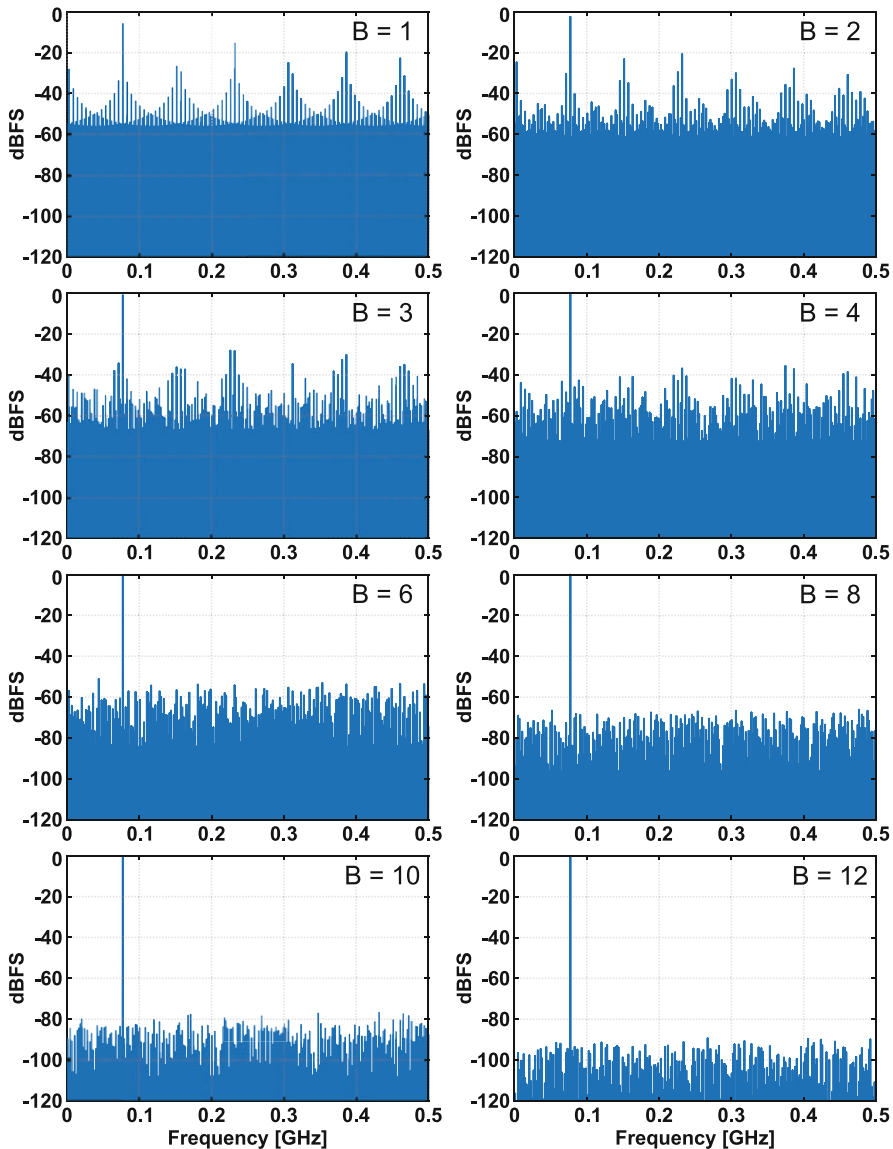


Fig. 2.11 Frequency spectra of an ideally quantized with various resolutions 77 MHz signal sampled at 1 GS/s ($N_{\text{FFT}} = 1024$)

spurs [19]. This is understood by the fact that for every added bit $\Delta^2/12$ reduces by 6 dB, while the additional 3 dB results from preserving the same total harmonic energy with twice the number of harmonics.

Table 2.1 Comparison between calculated and simulated SQNR for different B

Number of bits [B]	Calculated SQNR	Simulated SQNR
1	7.78 dB	6.31 dB
2	13.80 dB	13.30 dB
3	19.82 dB	19.53 dB
4	25.84 dB	25.61 dB
5	31.86 dB	31.66 dB
6	37.88 dB	37.73 dB
7	43.90 dB	43.84 dB
8	49.92 dB	49.84 dB
10	61.96 dB	61.92 dB
12	74.00 dB	73.98 dB

Having determined the conditions under which ϵ_q is considered white noise, the Signal-to-Quantization-Noise Ratio (SQNR) within the Nyquist bandwidth can be computed for a FS input sinusoid with a peak-to-peak amplitude of V_{FS}

$$\begin{aligned} SQNR &= 10 \log \left[\frac{(\frac{V_{FS}}{2\sqrt{2}})^2}{(\frac{V_{FS}}{\sqrt{12 \cdot 2^B}})^2} \right] = 10 \log(1.5 \cdot 2^{2B}) \quad [\text{dB}] \quad (2.9) \\ &= 6.02 \cdot B + 1.76. \end{aligned}$$

As anticipated, due to the non-linear nature of ϵ_q , the validity of the above expression may be questionable as the resolution decreases or for a signal that doesn't uniformly occupy a sufficient range [12]. Table 2.1 compares the calculated ideal SQNR against the simulated value for different resolutions. The noise approximation leading to Eq. (2.9) provides an overestimation, which reduces as the resolution increases, eventually matching the simulated value.

Finally, if the utilized signal bandwidth $f_{in,bw}$ does not include the complete Nyquist band, such that the sampling happens at a higher rate than Nyquist, there is an improvement in SQNR equivalent to the oversampling ratio $f_s/(2 \cdot f_{in,bw})$. In this case, an extra term known as the processing gain needs to be included in Eq. (2.9), which now becomes

$$SQNR = 6.02 \cdot B + 1.76 + 10 \log \left[\frac{f_s}{2 \cdot f_{in,bw}} \right]. \quad [\text{dB}] \quad (2.10)$$

Oversampling combined with quantization error shaping and digital filtering to remove out-of-band noise are fundamental concepts in $\Delta\Sigma$ converters [20].

2.2 Error Sources

Although ideally ϵ_q sets the theoretical single conversion error source, imperfections of electronic components utilized in a real A/D conversion introduce several noise and distortion sources to the signal. The sampling network comes with thermal noise, non-linear distortion, and aperture jitter. The actual quantizer introduces further thermal noise and both integral and differential non-linearity on top of its existing quantization noise. For very wide bandwidth, if an additional analog front-end needs to be utilized, it adds extra thermal noise and non-linear distortion. Figure 2.12 illustrates the model of a real converter including the aforementioned error sources.

2.2.1 Noise

The wideband internal circuits in a converter produce a certain amount of thermal noise due to Brownian motion of charges. Although the instantaneous value of noise cannot be predicted, its Gaussian nature allows for the construction of a statistical model by means of a distribution. To measure its RMS value a large number of output samples are collected and plotted as a histogram, from where the mean μ and the standard deviation σ (or variance σ^2) can be calculated² [21]. The RMS noise voltage is equal to σ and can be expressed either with respect to an LSB or as an RMS absolute voltage.

Three main noise sources can be identified in a converter chain (Fig. 2.12), namely, thermal noise from the sampling network; thermal noise due to the quantizer; and aperture jitter during the sampling instants.

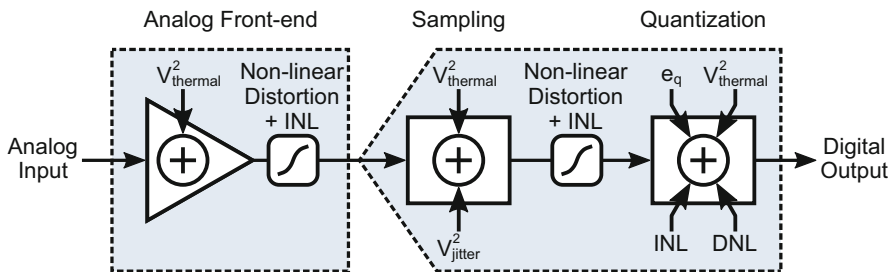


Fig. 2.12 Conceptual model of a real converter including error sources from the different blocks

² In the subsequent calculations, the noise variance will be expressed as voltage squared.

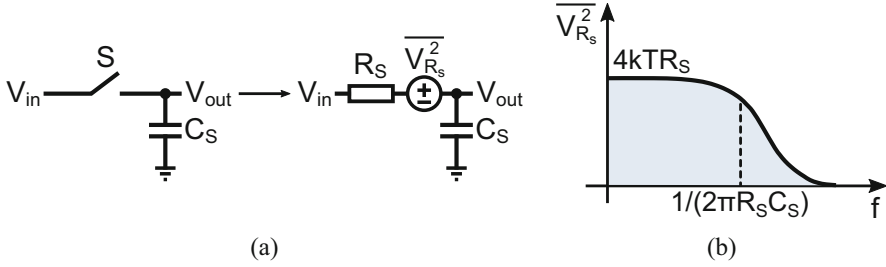


Fig. 2.13 (a) Simple model of a sampler and (b) its noise spectrum

Sampler Thermal Noise

The simplest implementation of a sampler comprises a switch S (Metal-Oxide-Semiconductor (MOS) device) and a capacitor C_S , as illustrated in Fig. 2.13a. When S is turned on, the MOS device is operating in triode region; therefore, it exhibits an on-resistance R_S . R_S produces white noise with a spectral density (single-sided) of

$$\overline{V_{R_S}^2} = 4kTR_S, \quad [\text{V}^2/\text{Hz}] \quad (2.11)$$

where $k = 1.38 \cdot 10^{-23}$ J/K is the Boltzmann constant and T is the absolute temperature.³ The RC network of the sampler shows a first-order low-pass characteristic with a cut-off frequency of

$$f_{-3\text{dB}} = \frac{1}{2\pi R_S C_S}, \quad [\text{Hz}] \quad (2.12)$$

which shapes the noise spectrum of R_S as shown in Fig. 2.13b. The sampler noise power can then be calculated by integrating $\overline{V_{R_S}^2}$ over the entire noise bandwidth

$$\overline{V_{n,\text{samp}}^2} = \alpha_{\text{FE}} \int_0^\infty \frac{4kTR_S}{(2\pi f R_S C_S)^2 + 1} df = \alpha_{\text{FE}} \frac{kT}{C_S}, \quad \alpha_{\text{FE}} \geq 1, [\text{V}^2] \quad (2.13)$$

where α_{FE} accounts for any excess noise in the presence of an analog front-end.

Quantizer Thermal Noise

A typical 1-bit quantizer employs a dynamic latch-based comparator (see Chap. 4) in some form and combination. To provide a simple expression as a basis for the noise of the quantizer, we construct the model shown in Fig. 2.14a. It assumes a two-stage comparator with a $g_{m,L}$ latch output and a $g_{m,I}$ integrator input [22] to provide some gain prior to regeneration and lower the noise of the latch. Ignoring large signal behavior and considering the latch as a settling stage with a $g_{m,L}$ noise

³ Throughout this book, T is set to 323 K (50 °C), unless explicitly stated otherwise.

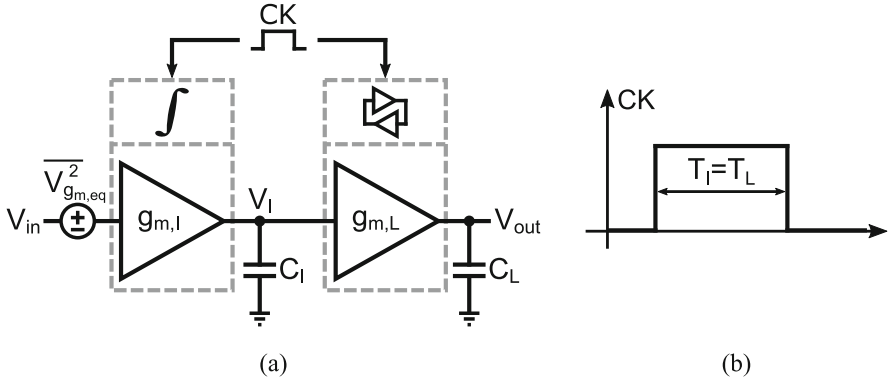


Fig. 2.14 (a) Simple quantizer model and (b) its allowed operation time

contribution equivalent to an effective resistor of $1/g_{m,L}$ [23], the latch noise power at V_1 is given by

$$\overline{V_{g_{m,L}}^2} = \frac{4kT\gamma}{g_{m,L}} \cdot \frac{g_{m,L}}{4C_L} = \frac{kT}{C_L}, \quad [\text{V}^2] \quad (2.14)$$

where γ is the thermal noise excess factor.⁴ The input stage integrates its own noise over a noise bandwidth proportional to $1/2T_1$, where T_1 is the integration time allowed for the quantizer (Fig. 2.14b). Its input noise power can be calculated similarly to [25] and is given by

$$\overline{V_{g_{m,I}}^2} = \frac{4kT}{g_{m,I}} \cdot \frac{1}{2T_1} = \kappa \frac{kT}{AC_1}, \quad [\text{V}^2] \quad (2.15)$$

where κ depends on the integration time, the integration voltage on V_1 , and the relative biasing of the input devices. Assuming for simplicity equal values for C_1 and C_L , the total input-referred noise power can be approximated as

$$\overline{V_{n,\text{quant}}^2} = \overline{V_{g_{m,I}}^2} + \frac{1}{A^2} \overline{V_{g_{m,L}}^2} \approx \frac{kT}{AC_1}, \quad [\text{V}^2] \quad (2.16)$$

where in the last step we substituted $\kappa = 1$ and $A = 4$ for the input stage.⁵

⁴ In literature, values of γ for short-channel devices span between 0.7 and 2.9 [24]. In this book, the value of 1 will be used unless otherwise stated.

⁵ The maximum gain for a g_m - C integrator cannot exceed the $g_m R_o$ of a differential pair, which in a 28 nm bulk Complementary Metal-Oxide-Semiconductor (CMOS) process can reach values of 4 (12 dB) at GHz operation.

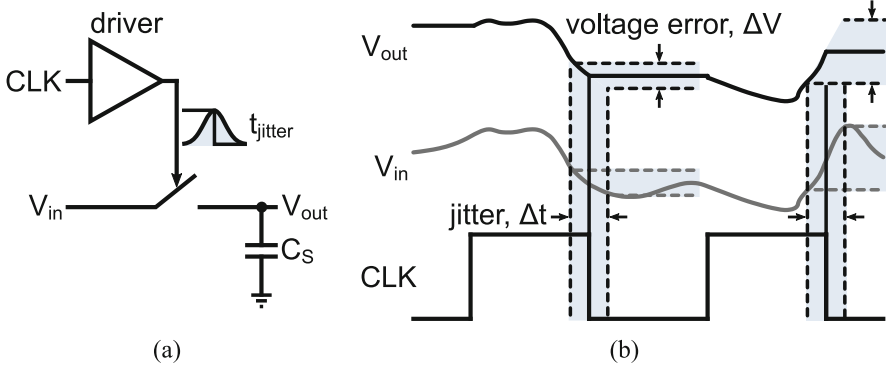


Fig. 2.15 (a) Sampler with jitter and (b) time to voltage error translation

Aperture Jitter

During ideal sampling (Sect. 2.1.1), the continuous-time input signal is sampled precisely at instants equally spaced by T_s . However, noise and mismatch in the devices of a real sampling network result in random variations in the clock edge (Fig. 2.15a), leading to sampling uncertainty, known as aperture uncertainty or aperture jitter. It is generally measured in picoseconds RMS. Jitter in time (Δt) translates into an output voltage error (ΔV), whose value is strongly related to the slope of the input signal, as illustrated in Fig. 2.15b. It is worth mentioning that jitter on the sampling clock or on the analog input produce exactly the same type of error. In fact, assuming that the sources are uncorrelated, they simply add in a Root-Sum-Square (RSS) fashion to yield the total error at the output.

The voltage error due to jitter can be easily calculated for a sinusoidal input of $V_{in}(t) = 0.5V_{FS}\sin(2\pi f_{in}t)$.⁶ Since this error depends on the slope of the signal, it is maximum at the zero crossings

$$\begin{aligned}\Delta V_{\max} &= \left. \frac{d}{dt} V_{in}(t) \cdot \Delta t \right|_{t=0} = 2\pi f_{in} \frac{V_{FS}}{2} \cos(2\pi f_{in}t) \cdot \Delta t \Big|_{t=0} \\ &= \pi f_{in} V_{FS} \cdot \Delta t.\end{aligned}\quad (2.17)$$

Since Δt is assumed to be random with a standard deviation of t_{jit} , the integrated error noise power can be approximated as

$$\begin{aligned}\overline{V_{n,jitter}^2} &= \frac{1}{T_{sig}} \int_0^{T_{sig}} \left(\frac{d}{dt} V_{in}(t) \right)^2 dt \cdot t_{jit}^2 \\ &= \frac{1}{2} (\pi f_{in} V_{FS})^2 \cdot t_{jit}^2,\end{aligned}\quad [V^2] \quad (2.18)$$

⁶ The calculation is done with respect to a peak-to-peak signal amplitude to preserve consistency with all our subsequent calculations.

where T_{sig} is the integration period, which for a sinusoid can be chosen as the signal period.

As a final note on jitter, special care must be taken across the entire input and clock chains to minimize the accumulative contribution of every added block. In Chap. 6, we will present an ultra-low jitter clock chain that shows how such a minimization can be achieved.

Now that we derived all the major noise contributions referred to the residue node (quantizer input), they can be summed and added to $\bar{\epsilon}_q$ to yield a first-order total quantization and noise power (single-ended)

$$\overline{V_{\epsilon_q+n,\text{total}}^2} = \frac{\Delta^2}{12} + \alpha_{\text{FE}} \frac{kT}{C_S} + \frac{kT}{AC_1} + \frac{1}{2} (\pi f_{\text{in}} V_{\text{FS}})^2 \cdot t_{\text{jitter}}^2 \cdot [V^2] \quad (2.19)$$

One quick observation arising from the above expression is that $\overline{V_{n,\text{jitter}}^2}$ increases with f_{in} , whereas both $\overline{V_{n,\text{samp}}^2}$ and $\overline{V_{n,\text{quant}}^2}$ are to a first-order input frequency independent. Additionally, to reduce both $\overline{V_{n,\text{samp}}^2}$ and $\overline{V_{n,\text{quant}}^2}$ the capacitors at the corresponding band-limiting nodes must increase, adversely affecting the bandwidth. Section 2.4 analyzes the accuracy degradation of a converter due to the above noise sources and establishes some fundamental accuracy-speed-power limits.

2.2.2 Non-linearity

The non-linearity of the circuit elements utilized in a real converter will make its transfer characteristic deviate from an ideal equal step width linear curve. As illustrated in Fig. 2.16, these deviations manifest themselves both locally in each step (Fig. 2.16a) and globally across the entire characteristic (Fig. 2.16b). The two main types of non-linearity encountered in a real converter are characterized by the Differential Non-Linearity (DNL) and the Integral Non-Linearity (INL)

DNL quantifies the individual deviation of each step's width from the ideal value Δ (1 LSB) according to the following expression:

$$DNL_i = \frac{(V_{i+1} - V_i) - \Delta}{\Delta}, \quad \forall i = 0 \dots (2^B - 2). \quad (2.20)$$

For each step, the relative deviation of its width from Δ is uncorrelated with the equivalent deviation of the previous and next steps. Positive or negative DNL implies a larger or smaller step compared to Δ , respectively. A value of -1 LSB is the smallest possible and indicates that a step was completely skipped, a situation described as a missing code (Fig. 2.16a). In the presence of a noisy signal, such that the transition levels carry noise comparable to Δ , this noise can affect the

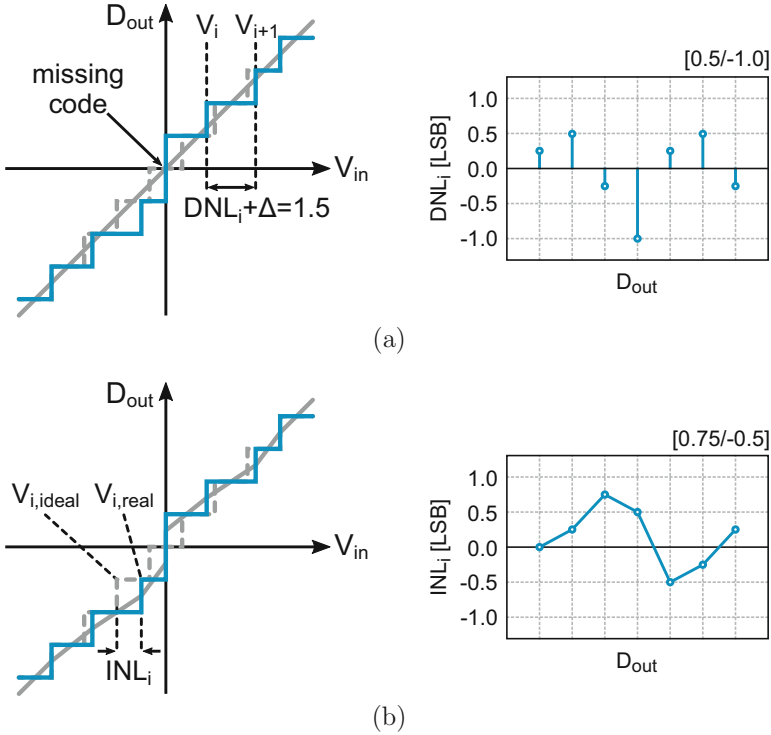


Fig. 2.16 (a) DNL in transfer characteristic with corresponding curve and (b) INL in transfer characteristic with corresponding curve

DNL true value and potentially hide missing codes [26, 27]. Therefore, its value alone should not be trusted blindly. DNL is due exclusively to the quantizer, and the ENC (Fig. 2.8) determines how its errors spread across the transfer curve. Strictly speaking, these errors result in distortion products at the converter output, which depend both on the amplitude of the signal and on their relative position along the transfer curve. However, similar to ϵ_q , under the assumption of a uniform DNL spread across the FS, its effect can be seen more as random noise rather than distortion. In that case, the degradation in SQNR can be estimated if a DNL within $\pm d$ is added to the signal, resulting in a worst-case total quantization + DNL error within $\pm 1/2(\Delta + d)$. Adding this to Eq. (2.9) results in the Signal-to-Quantization-and-DNL-Noise Ratio (SQDNR)

$$\begin{aligned}
 SQDNR &= 10 \log \left[\frac{\left(\frac{V_{FS}}{2\sqrt{2}}\right)^2}{\left(\frac{(1+d)V_{FS}}{\sqrt{12}\cdot 2^B}\right)^2} \right] = 10 \log \left(1.5 \cdot \frac{2^{2B}}{(1+d)^2} \right) [\text{dB}] \quad (2.21) \\
 &= 6.02 \cdot B - 1.76, \quad \text{if } d = 0.5.
 \end{aligned}$$

INL quantifies the overall deviation of the actual converter transfer characteristic from a straight line passing through the first and last transitions. Alternatively, if we draw a line passing through all the real transitions (Fig. 2.16b), its deviation from the ideal straight line (Fig. 2.16a) reveals INL. In each step, INL can be calculated as follows:

$$INL_i = \frac{(V_{i,real} - V_{i,ideal})}{\Delta}, \quad \forall i = 0 \dots (2^B - 1). \quad (2.22)$$

In contrast to the DNL, INL has a cumulative nature adding up errors from the consecutive steps to move the transfer curve with respect to the straight line, therefore resulting in an integral error. As such, its “purity” is affected less than the DNL in the presence of noise, making its value more trustworthy. It can be shown that INL from the quantizer only in each step can be calculated by a cumulative summation of the individual DNLs up to the previous step by the following expression:

$$INL_j = \sum_{i=0}^{j-1} DNL_i. \quad (2.23)$$

The total converter INL is a summation in RSS of different contributions from all the blocks in the chain that generate distortion, including the sampling network and the analog front-end (if utilized) (Fig. 2.12). It is not exclusively a quantizer property like DNL. Overall, INL results in input signal-dependent distortion products at the converter output, making it hard sometimes to identify which one of the individual contributors is dominant.

Non-monotonicity describes a special situation, where an increasing/decreasing input signal results in a decreasing/increasing step in the transfer curve, making the width of that step (hence its DNL) “ill-defined” [26]. This situation is especially important for converters used in closed-loop configurations; therefore, it should be avoided by design. It can be shown that a sufficient but not necessary condition for INL to prevent non-monotonicity is given below

$$|INL_i| \leq 0.5 \text{ LSB}, \quad \forall i, \quad (2.24)$$

which then results in an equivalent condition for DNL as follows:

$$|DNL_i| \leq 1 \text{ LSB}, \quad \forall i. \quad (2.25)$$

2.2.3 Calibration

Generally speaking, any type of non-linearity error, including DNL and INL, originate from circuit imperfections (mismatch [28], leakage, incomplete settling, voltage-temperature variations, etc.) and/or technology limitations to achieve a required performance. Their contribution can be minimized by proper design (e.g., device up-scaling) and/or architectural choices, which often increase the power consumption and area while compromising speed.

Alternatively, deterministic errors that are not associated with random noise but stem from circuit or technology imperfections may be compensated by means of calibration techniques. Such techniques can potentially yield a better overall performance with a reduced impact on the power consumption. The compensation process primarily comprises the following steps:

1. Error detection by measuring circuits' parameters that are considered for modification
2. Error correction by modifying the parameters to desired values by the correction circuitry, such that the errors are minimized or eliminated

The error detection can be implemented either in the analog or in the digital domain. The optimal implementation depends on the type and magnitude of errors as well as the application, performance, and technology at hand. Additional circuits and test signals are often necessary to perform the detection; however, it can be also performed by a statistical analysis without requiring extra hardware or modifications to the core circuitry. The error correction can be also performed either in the analog or in the digital domain (or a combination of both), with the two having distinct differences regarding the end result of the calibration and circuitry used. For example, if correction is performed in the analog domain, modifications in the core circuits are necessary in order to re-adjust the parameters (e.g., by changing biasing voltages/currents or adding/subtracting tunable loads) and eliminate the error. The loading effects of such modifications on the core circuits' performance must then be taken into account. If digital correction is performed, the core circuits are left untouched, and the inverse of the error function is digitally created and applied to the digital output to reduce the error. In this case, the calibration accuracy may be somewhat inferior due to rounding effects but with increasing power and speed benefits moving into finer CMOS processes.

A final difference lies with how often the calibration is performed and how disruptive it is to the normal operation. In case of the so-called "foreground" calibration, the converter operation is halted, and once the calibration is performed, it becomes available again to continue its operation. In the case of "background" calibration, the converter errors are corrected simultaneously to its normal operation, and the calibration is integrated ideally seamlessly into the core functionality. As expected, both methods have advantages and drawbacks in terms of hardware, signal range utilization, correction accuracy, and error tractability. Therefore, the optimal

choice depends on the nature of errors and the specific application requirements and tolerances.

2.3 Performance Evaluation

A converter's achievable performance can be evaluated in the time domain and in the frequency domain [3, 12, 15], and several metrics exist for such evaluations. Below, we will limit ourselves to the frequency domain evaluation by means of an FFT [29] and define the metrics that will be used in the following chapters.

2.3.1 Metrics

Nth-order Harmonic Distortion (HD_n) is normally specified in dBc (decibels below carrier) and is the reciprocal of the ratio between the RMS value of the fundamental signal and the RMS value of its nth-order harmonic. The harmonics of the input signal can be distinguished from other distortion products because of their location in the frequency spectrum at integer multiples of the input frequency. HD_n is generally specified for input signals near FS since for much smaller signals, there may be other error mechanisms that dominate.

Total Harmonic Distortion (THD) is the inverse ratio of the RMS value of the fundamental signal to the mean RSS value of its harmonics. Depending on the specific design and application, the first five to seven harmonics are considered significant. For a FS input sinusoid with a peak-to-peak amplitude of V_{FS} and harmonics' amplitude of $V_{\text{harm},n}$, $n = 2, 3, \dots, 7$, THD is evaluated by the following expression:

$$THD = -10 \log \left[\frac{\left(\frac{V_{FS}}{2\sqrt{2}}\right)^2}{\sqrt{V_{\text{harm},2}^2 + \dots + V_{\text{harm},7}^2}} \right]. \quad [\text{dB}] \quad (2.26)$$

Signal-to-Noise Ratio (SNR) is the ratio of the RMS signal amplitude to the mean RSS value of all noise-related spectral components, including quantization (plus DNL), thermal, and jitter. For a FS input sinusoid with a peak-to-peak amplitude of V_{FS} , its value is evaluated as

$$SNR = 10 \log \left[\frac{\left(\frac{V_{FS}}{2\sqrt{2}}\right)^2}{\sqrt{\epsilon_{q+DNL}^2 + V_{thermal}^2 + V_{jitter}^2}} \right]. \quad [\text{dB}] \quad (2.27)$$

Signal-to-Noise-and-Distortion Ratio (SNDR) or SINAD is the ratio of the RMS signal amplitude to the mean RSS value of all spectral components, including quantization error, noise, and harmonics. Again, for a FS input sinusoid with a peak-to-peak amplitude of V_{FS} , the following expression evaluates SNDR:

$$SNDR = 10 \log \left[\frac{\left(\frac{V_{FS}}{2\sqrt{2}}\right)^2}{\sqrt{V_{noise}^2 + V_{harmonics}^2}} \right]. \quad [\text{dB}] \quad (2.28)$$

There exists a relation between THD, SNR, and SNDR provided all of them are characterized under the same input signal conditions (amplitude and frequency) [30]. This relation is summarized with the equations below

$$THD = -10 \log \left[10^{-(SNDR/10)} - 10^{-(SNR/10)} \right], \quad [\text{dB}] \quad (2.29)$$

$$SNR = -10 \log \left[10^{-(SNDR/10)} - 10^{-(THD/10)} \right], \quad [\text{dB}] \quad (2.30)$$

$$SNDR = -10 \log \left[10^{-(SNR/10)} + 10^{-(THD/10)} \right]. \quad [\text{dB}] \quad (2.31)$$

Effective Number of Bits (ENOB) is the actual converter accuracy after adding up all error sources. It can be calculated by using Eq. (2.9) and solving for B after substituting SNDR for SQNR

$$ENOB = \frac{SNDR - 1.76}{6.02}. \quad (2.32)$$

Spurious Free Dynamic Range (SFDR) is one of the most important specifications in ADCs for communications applications. It is quantified as the ratio of the RMS value of the fundamental signal to the RMS value of the largest undesired spectral content. It may be specified either in dBc or in dBFS (decibels below FS). For input signals near FS, it typically coincides with the largest HDn. There might be cases though, where some other distortion product determines SFDR (e.g., an error tone due to interleaving; see Sect. 3.7 from the next chapter).

Analog Bandwidth (BW) is defined as the frequency at which the output power of the reconstructed fundamental drops by 3 dB below its low-frequency value. It does not contain any useful information regarding the spectral purity of the converter at that frequency.

Effective Resolution Bandwidth (ERBW) is defined as the frequency at which there is a 3 dB drop in SNDR (or a 0.5 bit drop in ENOB) compared to its low-frequency value. For reasons that will become obvious in the following chapter, it is highly desirable (but not always easily achievable) that both the analog BW and the ERBW are above the Nyquist frequency.

Noise Spectral Density (NSD) is another important frequency domain metric that measures the noise per unit bandwidth at a given frequency. It may be specified either in V^2/Hz or in dB/Hz . Assuming a flat NSD over a certain band, the SNR within this bandwidth is linked with the NSD via the expression

$$NSD = -SNR - 10 \log(BW). \quad [\text{dB}/\text{Hz}] \quad (2.33)$$

N^{th} -order Intermodulation Distortion (IM n) is the equivalent HD n when applying two closely spaced sinusoidal inputs at frequencies f_1 and f_2 . The amplitude of each tone is backed off by at least 6 dB compared to a one-tone to avoid clipping upon in-phase addition of the two tones. The second-order and third-order products are usually the dominant ones. The second-order products are located at $f_2 \pm f_1$ and can be removed by filtering. The third-order products contain two pairs located at $2f_1 \pm f_2$ and $2f_2 \pm f_1$, respectively. The ones at $2f_1 - f_2$ and $2f_2 - f_1$ are of special interest since they fall close to the two fundamentals and properly characterize the converter's spectral purity.

Multi-Tone Power Ratio (MTPR) can be seen as an evaluation metric for the in-band SFDR when multiple sinusoidal inputs are applied. This metric is particularly useful in multi-channel communication systems such as Orthogonal Frequency Division Multiplexing (OFDM) [31]. A large number of tones equal in amplitude and in frequency spacing are applied, and one of them is eliminated from the input signal leaving an empty bin [32]. However, due to the converter's distortion, a small signal appears in that bin. The ratio between the RMS value of one of the fundamental signals and the RMS value of the undesired spectral content in the empty bin yields the MTPR.

2.3.2 Figures of Merit

Some of the metrics described in Sect. 2.3.1 can be used in different combinations and ratios in order to compare the performance of different converters covering similar applications. For this reason, the Figure-of-Merit (FoM) concept has been introduced, serving to measure the power efficiency of a converter with respect

to other specifications, with speed (sample rate) and accuracy the dominant ones. Although many different FoMs exist, two are extensively used in literature and will be summarized below.

Walden's FoM Originally proposed in [33] for Nyquist converters and later adjusted to also cover oversampled converters [34], FoM_W is defined as

$$FoM_W = \frac{Power}{2^{ENOB} \cdot \min\{2BW, f_s\}} \quad [\text{J/conv.-step}] \quad (2.34)$$

and quantifies the energy spent by a converter to achieve a certain accuracy while performing the conversion at a certain speed. Its units are energy (in J) per conversion step. As Eq. (2.34) suggests, for every extra bit of ENOB, power increases by $2\times$. This trend is not obeyed by noise-limited converters, whose power would need to increase by $4\times$ (see Sect. 2.4), which is an important limitation of this FoM.

Schreier's FoM To alleviate the limitation regarding noise-limited converters, FoM_S was proposed, initially ignoring distortion [20] and later adjusted to include both noise and distortion [35]. It is defined as

$$FoM_S = SNDR + 10 \log \left[\frac{\min\{BW, f_s/2\}}{Power} \right]. \quad [\text{dB}] \quad (2.35)$$

Its units are accuracy (in dB) and it depicts more correctly the $4\times$ higher energy per 6 dB of SNDR increase, which is the prevailing trend in the highest-performance designs of recent years. An extensive ADC performance survey by gathering data from works published at the major scientific venues for more than 20 years has been carried out by Prof. Boris Murmann of Stanford University and can be found in [36].

2.4 Accuracy-Speed-Power Limits

In Sect. 2.3.2, it was argued that a converter's performance is a trade-off between accuracy, speed,⁷ and power. The key challenge lies in maximizing the product with accuracy and bandwidth in the numerator and power in the denominator or minimizing its reciprocal by simultaneously pushing all the three parameters as far as possible toward the desired directions.

⁷ It is assumed that for a certain sample rate (speed), the converter needs to achieve the required accuracy for a bandwidth of at least half of that sample rate, and this assumption is used in the equations and plots to follow.

$$\uparrow \left[\frac{\uparrow \text{Accuracy} \cdot \text{Speed} \uparrow}{\text{Power} \downarrow} \right] \iff \left[\frac{\text{Power} \downarrow}{\uparrow \text{Accuracy} \cdot \text{Speed} \uparrow} \right] \downarrow. \quad (2.36)$$

Several error sources were identified in Sect. 2.2, which degrade the accuracy of a real converter below the ideal quantization error. As discussed in the previous section, errors that are associated with mismatch⁸ or non-linearity can be compensated either by design or by calibration with a small overhead on the other two parameters. On the other hand, errors stemming from thermal noise introduce a more fundamental trade-off on Eq. (2.36); improving one of its parameters will most likely result in an analogous degradation of the other two. The significance of such errors on the accuracy-speed-power are analyzed, and some fundamental limits on a converter's performance are established.

2.4.1 Sampler Noise Limit

In Sect. 2.2, Eq. (2.13) was derived for the single-ended sampler thermal noise. We repeat this expression here for a differential configuration,⁹ which is the start for our derivations, assuming an ideal noiseless front-end ($\alpha_{FE} = 1$)

$$\overline{V_{n,\text{samp}}^2} = \frac{2kT}{C_S}, \quad [\text{V}^2] \quad (2.37)$$

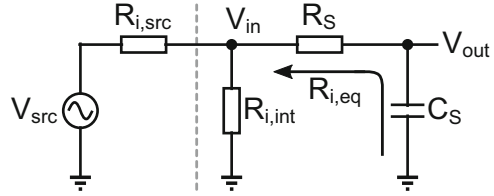
The accuracy degradation due to $\overline{V_{n,\text{samp}}^2}$ can be calculated by combining Eqs. 2.27 and 2.32 and considering a differential peak-to-peak signal swing of $V_{\text{FS-diff}}$

$$\begin{aligned} ENOB_{n,\text{samp}} &= \frac{1}{6.02} \cdot \left[10 \log \left(\frac{1}{8} \frac{V_{\text{FS-diff}}^2}{\epsilon_q^2 + \overline{V_{n,\text{samp}}^2}} \right) - 1.76 \right] \\ &= \frac{1}{6.02} \cdot \left[10 \log \left(\frac{1}{8} \frac{V_{\text{FS-diff}}^2}{\epsilon_q^2} \cdot \frac{1}{1 + \frac{\overline{V_{n,\text{samp}}^2}}{\epsilon_q^2}} \right) - 1.76 \right] \quad (2.38) \\ &= B - \frac{1}{6.02} \cdot 10 \log \left(1 + \frac{24}{C_S} \frac{kT}{2^{2B}} \right). \end{aligned}$$

⁸ A comprehensive analysis on the implications of mismatch in the design of analog circuits can be found in [37].

⁹ For differential signaling, the signal power increases by 4x, while the noise increases by 2x, leading to a 3 dB SNR improvement. Furthermore, the even-order harmonics are ideally fully suppressed, leading to an SFDR boost. On the downside, the power increases by 2x.

Fig. 2.17 Simple sampler model with input termination network



The minimum capacitance for a tolerable ENOB reduction can then be obtained for a certain input swing. It is evident from the above expression that to minimize the accuracy degradation due to $\overline{V_{n,samp}^2}$, C_S must be maximized. On the other hand, Eq. (2.12) implies that in order to maximize the bandwidth, C_S must be minimized (for a fixed R_S). To quantify this fundamental trade-off more completely, we add in the simple sampler model (Fig. 2.13) the basic input termination network, as shown in Fig. 2.17, which in some form is given in every converter measurement system. C_S can then be written as

$$C_S = \frac{1}{2\pi[(R_{i,src}/R_{i,int}) + R_S]f_{in}} = \frac{1}{\pi(0.5R_{i,int} + R_S)f_s}, \quad [F] \quad (2.39)$$

where $R_{i,src} = R_{i,int}$ and represent the external source resistance and the internal termination, respectively. Employing Eq. (2.28) with $\overline{V_{n,samp}^2}$ the sole noise contribution, and combining Eqs. (2.37) and (2.39), we reach to the final accuracy-speed limit

$$SNDR_{n,samp} = 10 \log \left[\frac{V_{FS-diff}^2}{8\pi kT(0.5R_{i,int} + R_S)f_s} \right]. \quad [dB] \quad (2.40)$$

The outcome of the above expression is that for a fixed termination network and C_S value, the only optimization “knob” in preserving the Nyquist $SNDR_{samp}$ as the sample rate increases is to reduce R_S accordingly. In Chap. 5, a sampling circuit that outperforms existing circuits in minimizing R_S will be presented.

The absolute minimum power required to charge C_S can be calculated in a similar fashion as in [38]. We assume that the charging occurs within half a period of f_s and the signal utilizes an input swing V_{FS} equal to the supply voltage V_{DD} . Keeping the $SNDR_{samp}$ as a measure of accuracy, the minimum power to achieve a certain accuracy dictated by the sampler noise is given by

$$\begin{aligned} P_{n,samp} &= V_{DD} \cdot I_{samp} = 2 \cdot 8 \cdot \overline{V_{FS}^2} \cdot f_s \cdot C_S \\ &= 16kTf_s \cdot SNDR_{n,samp}, \end{aligned} \quad [W] \quad (2.41)$$

where we substitute $SNDR_{n,samp} = \overline{V_{FS}^2}/\overline{V_{n,samp}^2}$. The above expression gives the accuracy-power limit due to the sampler noise. We can obtain the same result by allocating a full quantization noise contribution to the sampler and substituting C_S

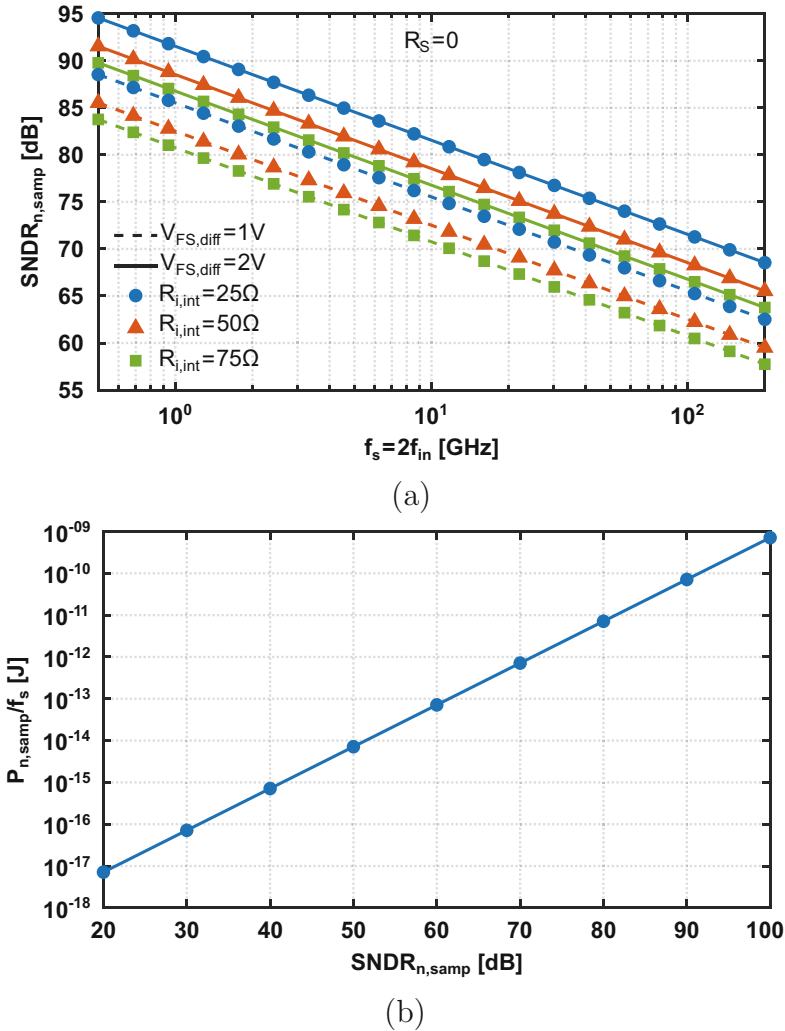


Fig. 2.18 Fundamental limits due to sampler noise: (a) accuracy-speed and (b) accuracy-power

in the above expression. The fundamental limits described by Eqs. 2.40 and 2.41 are plotted in Fig. 2.18 sweeping different parameters.

It is worth mentioning that recently published works [39–41] have shown progress in attempting to “break” the $V_{n,samp}^2$ fundamental limits described above. The underlying principle is to either decouple the generating noise source from the sampling bandwidth or sample the noise and then somehow cancel it. As such, these techniques necessitate additional components (resistors, capacitors, switches, amplifiers) in either open-loop or closed-loop configurations. When going at very high sample rates (>GHz), achieving the necessary amplification

and/or generating extra clocks (including associated routing overhead) for complex switching schemes, to bring down the noise, might take away some or all of the power, bandwidth, and area benefits of scaling down C_S . These might explain why such designs have yet to achieve sample rates beyond several MS/s.

2.4.2 Quantizer Noise Limit

The quantizer thermal noise introduces a second fundamental converter accuracy-speed-power limit. It is mainly defined by the input integrator stage preceding the final latch, as we also derived for our simple model of Fig. 2.14. This also makes the quantizer analysis easier, separating the noise critical input from the bandwidth critical latch (see Sect. 2.4.3). The two stages will be analyzed separately as they both impose different limits, and their contributions will be quantified. The noise power with all the assumptions from our basic model is written here in its differential form to start our derivations and given by

$$\overline{V_{n,\text{quant}}^2} = \frac{2kT}{AC_I}. \quad [\text{V}^2] \quad (2.42)$$

The accuracy reduction due to $\overline{V_{n,\text{quant}}^2}$ can be calculated by combining Eqs. 2.27 and 2.32 and considering a differential peak-to-peak signal swing of $V_{\text{FS-diff}}$

$$\begin{aligned} ENOB_{n,\text{quant}} &= \frac{1}{6.02} \cdot \left[10 \log \left(\frac{1}{8} \frac{V_{\text{FS-diff}}^2}{\epsilon_q^2 + \overline{V_{n,\text{quant}}^2}} \right) - 1.76 \right] \\ &= \frac{1}{6.02} \cdot \left[10 \log \left(\frac{1}{8} \frac{V_{\text{FS-diff}}^2}{\epsilon_q^2} \cdot \frac{1}{1 + \frac{\overline{V_{n,\text{quant}}^2}}{\epsilon_q^2}} \right) - 1.76 \right] \\ &= B - \frac{1}{6.02} \cdot 10 \log \left(1 + \frac{24 \frac{kT}{C_I}}{\frac{V_{\text{FS-diff}}^2}{2^B}} \right), \end{aligned} \quad (2.43)$$

which yields the minimum capacitance at the integrator output for a targeted reduction in ENOB and a given signal swing. To minimize this reduction, C_I ¹⁰ must be maximized, which adversely affects the input integrator's operating frequency, expressed as

¹⁰ Our model assumed $C_I = C_L$, which is not far from a realistic design scenario in 28 nm CMOS (see Chap. 4).

$$f_I = \frac{I_I}{C_I \Delta V_I} = \frac{g_{m,I} V_{GT,I}}{2C_I \Delta V_I}. \quad [\text{Hz}] \quad (2.44)$$

ΔV_I is the common-mode voltage rise/fall at the integrator output to build a certain gain, and I_I follows the basic MOS equation [42]

$$\frac{g_m}{I_D} = \frac{2}{V_{GT}}, \quad V_{GT} = \begin{cases} 2nkT/q \approx 60 - 80 \text{ mV}, & \text{Weak - Inversion} \\ V_{GS} - V_{TH}, & \text{Strong - Inversion} \\ 2(V_{GS} - V_{TH}), & \text{Velocity - Saturation} \end{cases}. \quad (2.45)$$

As with the sampler, we allocate half a period of f_s to the quantizer; thus, this is the maximum available time for the integrator. Combining Eqs. (2.42) and (2.44) and employing Eq. (2.28) with $\overline{V_{n,\text{quant}}^2}$ its only noise contribution, the accuracy-speed limit is derived

$$SNDR_{n,\text{quant}} = 10 \log \left[\frac{g_{m,I} V_{GT,I} V_{FS,\text{diff}}^2}{32kT \Delta V_I f_s} \right]. \quad [\text{dB}] \quad (2.46)$$

The minimum necessary power to charge C_I can be calculated with a similar method as for Eq. (2.41), following the same assumptions about the input signal. Additionally, by allocating a maximum value of $\Delta^2/12$ to $\overline{V_{n,\text{quant}}^2}$ for convenience,¹¹ the minimum power to achieve a certain accuracy dictated by the quantizer noise (accuracy-power limit) can be found as

$$\begin{aligned} P_{n,\text{quant}} &= V_{DD} \cdot I_I = 2 \cdot V_{FS} \cdot f_s \cdot C_I \cdot \Delta V_I \\ &= 4 \cdot V_{FS} \cdot f_s \cdot \frac{12kT}{V_{FS}^2} \cdot 2^{2ENOB_{n,\text{quant}}} \cdot \frac{V_{FS}}{2} \\ &= 24kT f_s \cdot 2^{\frac{SNDR_{n,\text{quant}} - 1.76}{6.02}}, \end{aligned} \quad [\text{W}] \quad (2.47)$$

where Eq. (2.32) is used, V_{DD} is assumed to be equal to V_{FS} , and ΔV_I is assumed to be half V_{FS} at the end of the integration. The fundamental limits described by Eqs. (2.46) and (2.47) are plotted in Fig. 2.19 sweeping different parameters. In Sect. 2.4.7, all limits will be plotted together for comparison.

2.4.3 Metastability Limit

In addition to the noise, metastability is another fundamental error source associated with the output latch stage of the quantizer. The latch regenerates exponentially on

¹¹ In high-speed converters, it is general practice to design the various thermal noise sources in the same order as the quantization noise.

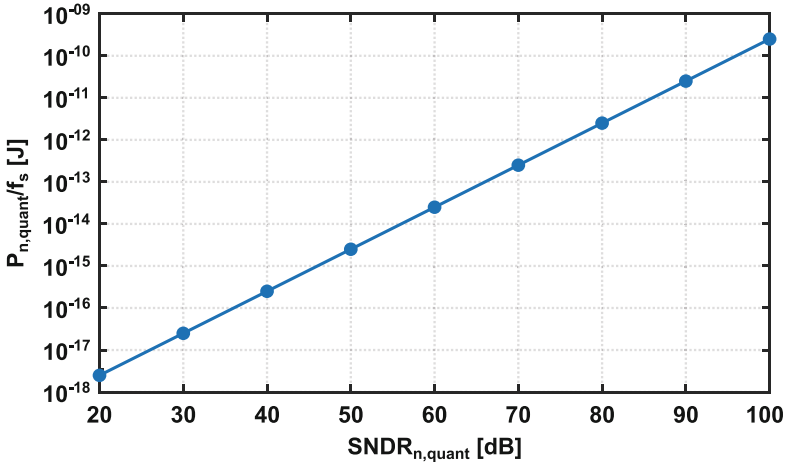
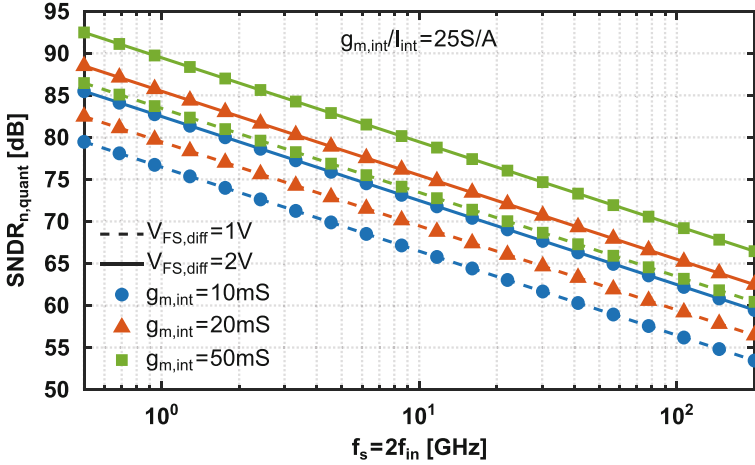


Fig. 2.19 Fundamental limits due to quantizer noise: (a) accuracy-speed and (b) accuracy-power

an input according to the following expression:

$$V_{out} = A V_{in} \cdot e^{-\frac{T_L}{\tau}} = A V_{in} \cdot e^{-\frac{s_{m,L} T_L}{C_L}}, \quad [V] \quad (2.48)$$

where A is the integrator's gain (see Sect. 2.2.1), while the time constant $\tau = C_L/g_{m,L}$ is a measure of the latch's bandwidth. Metastability refers to the situation where the quantizer differential input is so small (e.g., a fraction of an LSB), such that for the allowed operation time, the latch of the quantizer cannot produce a sufficiently large differential output for the following circuitry to unambiguously

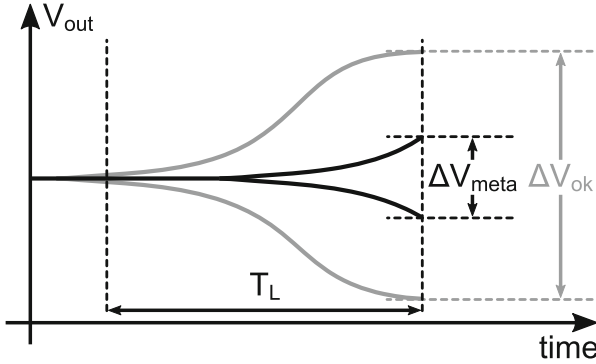


Fig. 2.20 Quantizer output for a valid (gray) and a metastable (black) case

perceive it as a clear logical level. This scenario, portrayed in Fig. 2.20, results in a conversion error, therefore leading to accuracy degradation. For a certain input voltage and a fixed gain A , this error can be reduced either by allowing more time to the quantizer to produce a sufficiently large output difference or by minimizing τ .

The error due to metastability may be interpreted as an increased quantization noise floor with a variance $\overline{\epsilon_q^2}$ multiplied by a certain probability of occurrence $PR(meta)$ [43]. The total converter noise may be then written as

$$\overline{V_{q+meta}^2} = \overline{\epsilon_q^2} \cdot [1 + PR(meta)]. \quad [V^2] \quad (2.49)$$

The second term inside the square brackets denotes the excess noise due to metastability. If we consider a differential input signal uniformly distributed within $\pm V_{FS-diff}/2$, then the probability of a metastable occurrence, otherwise known as Bit Error Rate (BER), can be seen as the ratio of the smallest input the latch can correctly regenerate on its given time divided by the full input range. For a B -bit quantizer with an equal probability of showing metastability in any of the 2^B steps, utilizing Eq. (2.48), $PR(meta)$ can be expressed as

$$PR(meta) = BER \cdot 2^{B_{meta}} = \frac{2^{B_{meta}} V_{in,min}}{\frac{V_{FS-diff}}{2^{B_{meta}+1}}} = \frac{2^{2B_{meta}} \cdot e^{-T_L/\tau}}{A}, \quad (2.50)$$

where it is assumed that the quantizer latch regenerates to V_{FS} and $B = B_{meta}$. $PR(meta)$ has an exponential dependency on τ ; therefore, minimizing it is extremely desirable. Further, if we re-write τ lumping the total capacitance at the quantizer output, we can see that the technology ultimately dictates the minimum achievable value

$$\tau \approx \frac{C_{gg}}{g_{m,L}} \approx \frac{1}{2\pi f_T}, \quad (2.51)$$

where f_T is the cut-off frequency for which the current gain is unity. In order to take into account practical limitations (e.g., layout parasitics), a more realistic value of $1/\pi f_T$ is adopted for τ , in all the subsequent analysis. Substituting Eqs. (2.50) and (2.51) into (2.49), we have

$$\overline{V_{q+meta}^2} = \overline{\epsilon_q^2} \cdot \left[1 + \frac{2^{2B_{meta}} \cdot e^{-\pi f_T / 2 f_s}}{A} \right], \quad [V^2] \quad (2.52)$$

where half a period of f_s is allocated for latch regeneration.¹² By allocating a certain small LSB fraction $a_{er} < 1$ to the error due to the excess noise in the above expression, and employing Eq. (2.32), the accuracy-speed limit imposed by metastability can be derived for various f_T values

$$SNDR_{meta} = 6.02 \left[\frac{\log_2(a_{er} A)}{2} + \frac{\pi f_T}{4 f_s \ln 2} \right] + 1.76. \quad [dB] \quad (2.53)$$

The take from the above expression is that if the quantizer resolution increases while preserving the same f_s and f_T , there is an increased excess noise due to metastability on the total quantization noise.

It is important to clarify that the above limit is derived under the assumption of f_s being the sample rate of a standalone non-pipelined non-interleaved quantizer. As such, it is the reciprocal of the standalone quantizer's latch delay T_L to achieve a certain resolution. Pipelining can improve this limit by reducing the quantizer resolution per pipeline stage, therefore increasing the overall resolution for the same total f_s or increasing the total f_s for the same overall resolution. Interleaving can also improve this limit, as discussed in the next chapter. By multiplexing several quantizers in time, each running at a lower standalone f_s , the aggregate f_s can be increased by the interleaving factor while also preserving the resolution.

In order to estimate the minimum power required by the latch to resolve within half a period of f_s a certain small input $A V_{FS-diff} / 2^{B_{meta}+1}$ ($A = 4 = 2^2$) and regenerate to V_{FS} , we start the derivation by substituting this value in Eq. (2.48) and solve for $g_{m,L}$

$$g_{m,L} = 2(B_{meta} - 2) \cdot \ln 2 \cdot f_s \cdot C_L. \quad [S] \quad (2.54)$$

This $g_{m,L}$ will require a minimum current I_L , and these two are related through the basic MOS Eq. (2.45). Before we reach to the final expression for the power, we need to substitute C_L from the latch noise Eq. (2.14) and assume that the input-referred latch noise voltage is at least $4\times$ smaller than the input that leads to metastability. This assumption aligns well with our two-stage quantizer model

¹² In a practical design, the input stage and the latch will each occupy a portion of the quantizer allocated time. Our simplification will affect our derivations by about $2\times$, which is tolerable for first-order generic derivations.

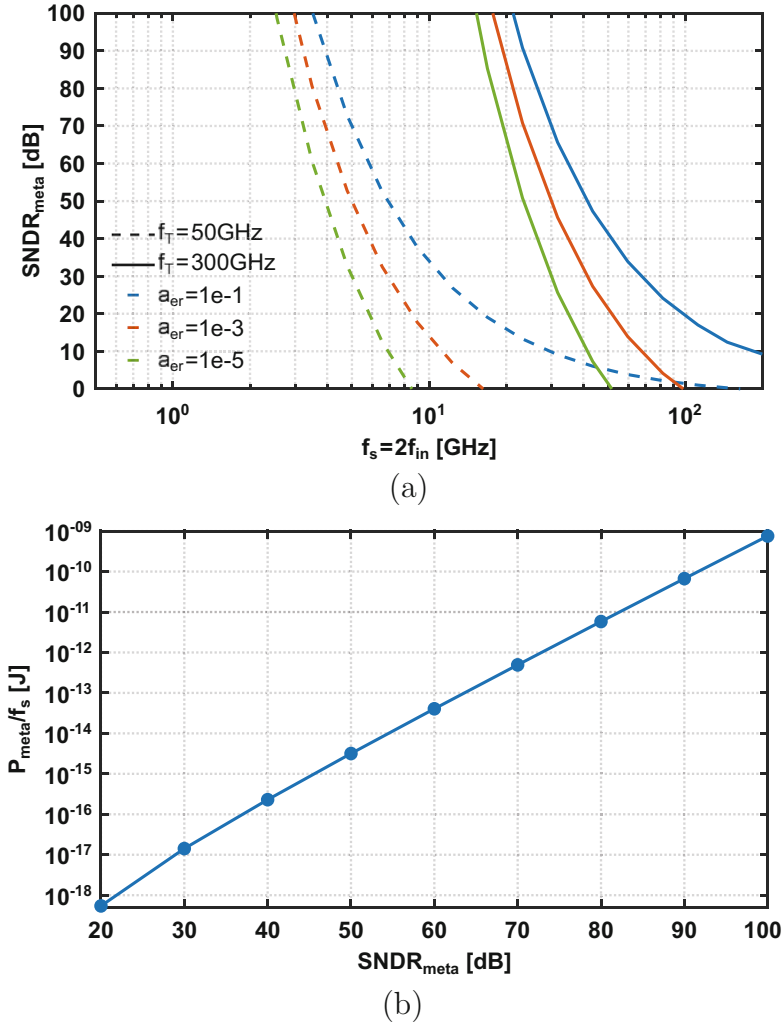


Fig. 2.21 Fundamental limits due to metastability of a standalone quantizer: (a) accuracy-speed and (b) accuracy-power

and allows to a first-order a proper metastability assessment. Finally, utilizing a supply voltage $V_{DD} = 1$ V equal to V_{FS} , we obtain the minimum power dictated by metastability, translating to the accuracy-power limit

$$\begin{aligned}
 P_{\text{meta}} &= V_{DD} \cdot I_L \\
 &= 24(B_{\text{meta}} - 2) \cdot \ln 2 \cdot f_s \cdot kT \cdot \frac{V_{GT}}{V_{FS}} \cdot 2^{2B_{\text{meta}}}. \quad [\text{W}] \quad (2.55)
 \end{aligned}$$

The fundamental metastability limits described by Eqs. (2.53) and (2.55) are plotted in Fig. 2.21 for different values of f_T and a_{er} .

2.4.4 Aperture Jitter Limit

Equation (2.18), from which the error noise power for a sinusoidal signal was obtained, can be adjusted to yield the differential jitter noise power for a differential peak-to-peak signal swing of $V_{\text{FS-diff}}$

$$\overline{V_{\text{n,jitter}}^2} = \frac{1}{2} (\pi f_{\text{in}} V_{\text{FS-diff}})^2 \cdot t_{\text{jit}}^2. \quad [\text{V}^2] \quad (2.56)$$

The accuracy reduction due to $\overline{V_{\text{n,jitter}}^2}$ can be calculated in a similar way as in Eqs. (2.38) and (2.43)

$$\begin{aligned} ENOB_{\text{n,jitter}} &= \frac{1}{6.02} \cdot \left[10 \log \left(\frac{1}{8} \frac{V_{\text{FS,diff}}^2}{\epsilon_{\text{q}}^2 + \overline{V_{\text{n,jitter}}^2}} \right) - 1.76 \right] \\ &= B - \frac{1}{6.02} \cdot 10 \log \left(1 + 2^{2B} \cdot 6(\pi f_{\text{in}})^2 \cdot t_{\text{jit}}^2 \right), \end{aligned} \quad (2.57)$$

from which the jitter value is obtained for a tolerable ENOB degradation and at a certain input frequency. The voltage error due to jitter is an increasing function of the frequency. This can be intuitively understood by the fact that a fixed error in time results in a larger voltage error when reflected to a signal with a faster slope compared to a slower slope signal. If we substitute Eq. (2.56) in the SNDR expression (Eq. (2.28)) and consider $\overline{V_{\text{n,jitter}}^2}$ the only noise source, the accuracy-speed limit due to jitter can be obtained

$$SNDR_{\text{n,jitter}} = 10 \log \left[\frac{1}{4(\pi f_{\text{in}})^2 \cdot t_{\text{jit}}^2} \right], \quad [\text{dB}] \quad (2.58)$$

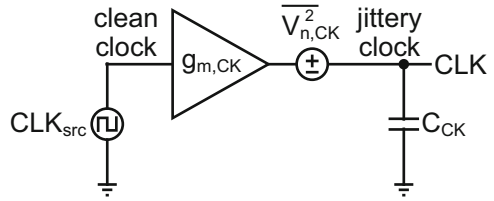
which is an already known expression [26], re-verified here by our analysis.

The minimum power to achieve a certain accuracy imposed by jitter noise is not entirely straightforward because strictly speaking, this power is not dissipated in the core converter parts (sampler and quantizer) but in the clock generation. Nevertheless, since the clock is an imperative part in any converter,¹³ we are including it in our fundamental limits for a comparison point.

To provide a first-order estimation of the clock power for a certain jitter, we model the clock generation as a single $g_{\text{m,CK}}$ unity gain buffer (Fig. 2.22) and assume linear operation for the entire clock swing, which is equal to V_{DD} . To simplify the analysis, we also assume that the dominant source leading to jitter is the

¹³ In the case of continuous-time ADCs [44, 45], although the input is not sampled, sampling is still performed along the chain (quantizer, back-end, reconstruction filter) to align with a synchronous clock. Depending on the part of the chain, jitter requirements can be different.

Fig. 2.22 Simple model for clock power estimation for a certain jitter



buffer thermal noise $\overline{V_{n,CK}^2}$, which, due to the unity gain, can be directly referred to the output. This noise can be calculated in a similar way as the quantizer noise (see Sect. 2.2.1, Eq. (2.16)). The buffer needs to charge C_{CK} ¹⁴ to V_{DD} , and we allocate a maximum of a quarter period of f_s to allow sufficient time for the actual sampling within half a period of f_s . The minimum required power consumed in the clock is then given as

$$\begin{aligned} P_{\text{jitter}} &= V_{DD} \cdot I_{CK} = 4V_{DD}^2 \cdot f_s \cdot C_{CK} \\ &= 4V_{DD}^2 \cdot f_s \cdot \frac{kT \cdot T_s^2}{16V_{DD}^2 \cdot t_{\text{jit}}^2}, \end{aligned} \quad [\text{V}^2] \quad (2.59)$$

where the Slew Rate (SR), which translates $\overline{V_{n,CK}^2}$ to t_{jit}^2 , has been written as voltage/time to provide V_{DD} within $0.25T_s$. By substituting t_{jitter}^2 from Eq. (2.56) for an input swing V_{FS} equal to V_{DD} and a Nyquist input frequency, the minimum power for a certain jitter is obtained

$$P_{\text{jitter}} = \frac{\pi^2}{2} kT f_s \cdot SNDR_{n,\text{jitter}}, \quad [\text{V}^2] \quad (2.60)$$

where $SNDR_{n,\text{jitter}} = \overline{V_{FS}^2} / \overline{V_{n,\text{jitter}}^2}$. Despite the several assumptions made to simplify the analysis, the above expression yields to a first-order a correct accuracy-power limit due to jitter, which is on par with the equivalent limits from the sampler and quantizer noise. The jitter-imposed limits of Eqs. (2.58) and (2.60) are plotted in Fig. 2.23 for several different parameters.

2.4.5 Mismatch Limit

At the beginning of this section, it was argued that errors associated with mismatch can be compensated with a small overhead, thus not introducing a fundamental trade-off between accuracy, speed, and power. Nevertheless, it is insightful to quantify the accuracy-speed and accuracy-power limits imposed by mismatch and

¹⁴ This capacitor includes the intrinsic buffer load and the sampling switch gate load.

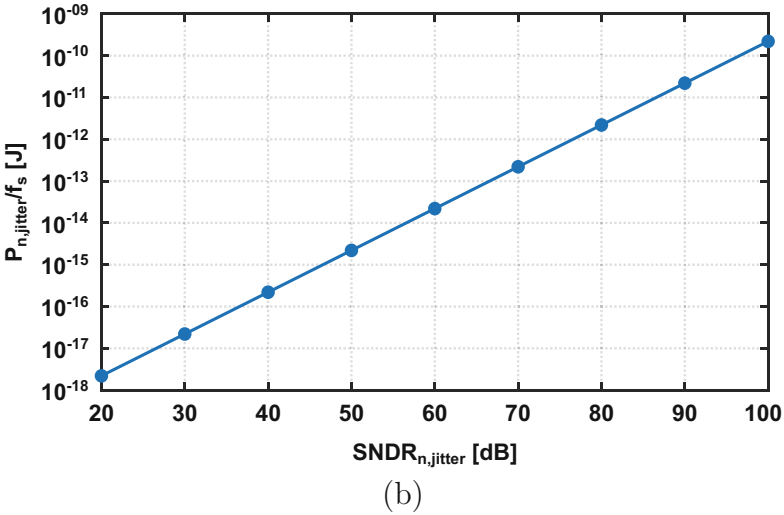
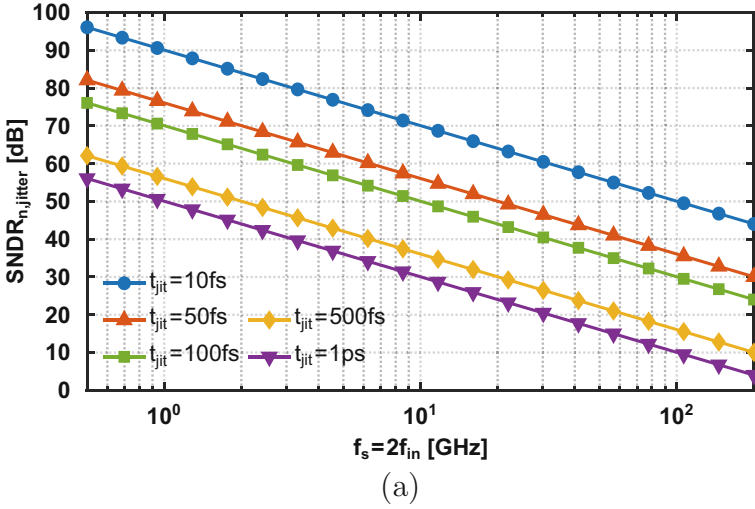


Fig. 2.23 Fundamental limits due to aperture jitter: (a) accuracy-speed and (b) accuracy-power

compare them to the derived ones imposed by noise, especially since the former are process dependent.

Similar to noise, mismatch is a random process as well, with a mean μ_M and a standard deviation σ_M (or variance σ_M^2). Assuming a differential pair with a mismatch dominated by the random variation in V_{TH} between the two devices, from Pelgrom's law [28], we obtain the variance

$$\sigma_M^2 = \frac{A^2 V_{TH}^2}{WL}, \quad [V^2] \quad (2.61)$$

where $A_{V_{TH}}$ is a mismatch constant that depends on the process. σ_M^2 is inversely proportional to the area. Assuming also that the devices of the differential pair are biased in strong inversion, the input capacitance C_M is found

$$C_M = (2/3)WLC_{ox} = \frac{2A_{V_{TH}}^2 C_{ox}}{3\sigma_M^2}. \quad [F] \quad (2.62)$$

This capacitance together with the source and internal termination resistances creates an upper limit to the input bandwidth, as shown in Eq.(2.39) if C_S is replaced by C_M . If we then combine Eqs.(2.28), (2.39), and (2.62) and consider a $3\sigma_M$ confidence interval for the mismatch contribution, we finally reach to the accuracy-speed limit

$$SNDR_{\sigma,match} = 10 \log \left[\frac{V_{FS-diff}^2}{48\pi A_{V_{TH}}^2 C_{ox} (0.5R_{i,int} + R_S) f_s} \right]. \quad [dB] \quad (2.63)$$

The minimum power required to charge C_M can be derived similarly to the one for charging C_S in the sampler noise limit. We allocate half a period of f_s for the operation and assume an input swing V_{FS} equal to the supply voltage. If we also consider a $3\sigma_M$ mismatch confidence interval, re-employing Eq.(2.41) and keeping $SNDR_{\sigma,match}$ as a measure of accuracy, we end up with the accuracy-power limit due to mismatch

$$\begin{aligned} P_{\sigma,match} &= V_{DD} \cdot I_{match} = 2 \cdot 8 \cdot \overline{V_{FS}}^2 \cdot f_s \cdot C_M \\ &= 48A_{V_{TH}}^2 C_{ox} f_s \cdot SNDR_{\sigma,match}, \end{aligned} \quad [W] \quad (2.64)$$

Comparing the above two expressions with Eqs.(2.40) and (2.41) giving the equivalent limits due to noise, we see $A_{V_{TH}}^2 C_{ox}$ in the denominator instead of kT , plus an extra multiplication factor depending on the targeted σ_M confidence interval. Both $A_{V_{TH}}$ and C_{ox} are technology-dependent parameters, indicating the effect of the process on the matching limit, in contrast to the baseline noise limit. Table 2.2 shows typical values of these parameters for three different process nodes [12], while the derived mismatch limits are plotted in Fig. 2.24.

Table 2.2 Typical process parameters and comparison with kT

Process [nm]	$A_{V_{TH}}$ [mV- μ m]	C_{ox} [fF/ μ m ²]	$A_{V_{TH}}^2 C_{ox}/kT$
130	5	11	61.7
65	4	13	46.7
28	2	25	22.4

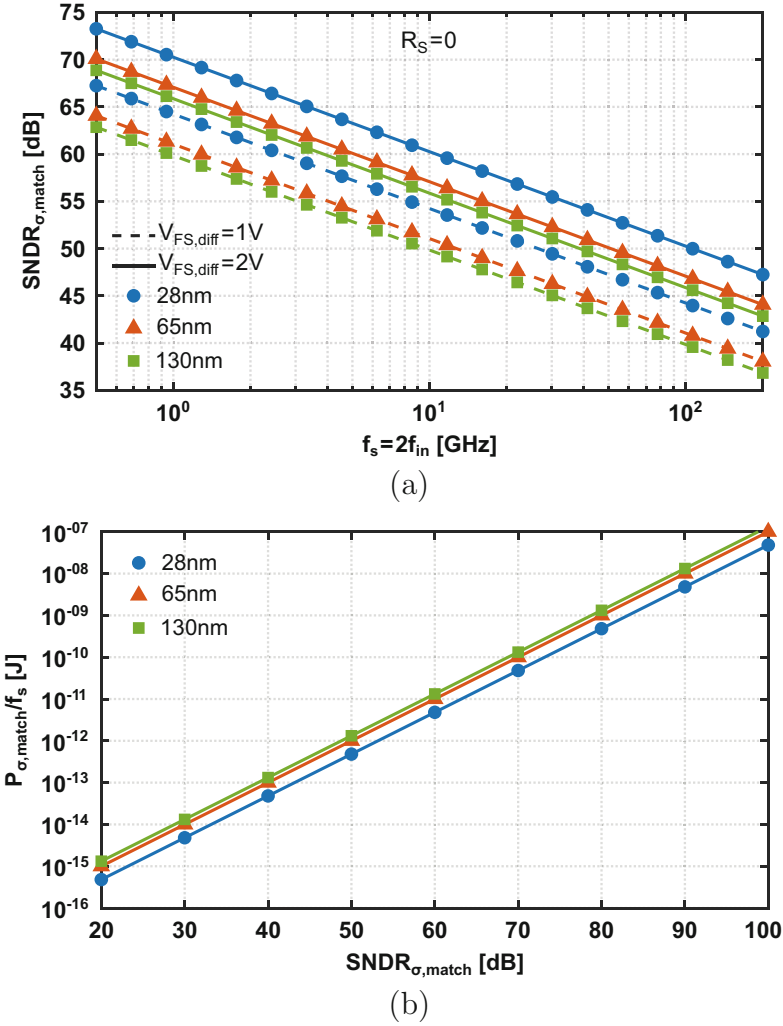


Fig. 2.24 Limits imposed by mismatch: (a) accuracy-speed and (b) accuracy-power

2.4.6 Heisenberg Uncertainty Principle

To complete our analysis, the Heisenberg uncertainty principle is also discussed based on [33], as the ultimate accuracy-speed limit in a converter’s performance, ultimately imposed by physics. The original principle [46] limiting what can be simultaneously known about the position and momentum of a quantum particle also applies to the energy-time complementary set stating

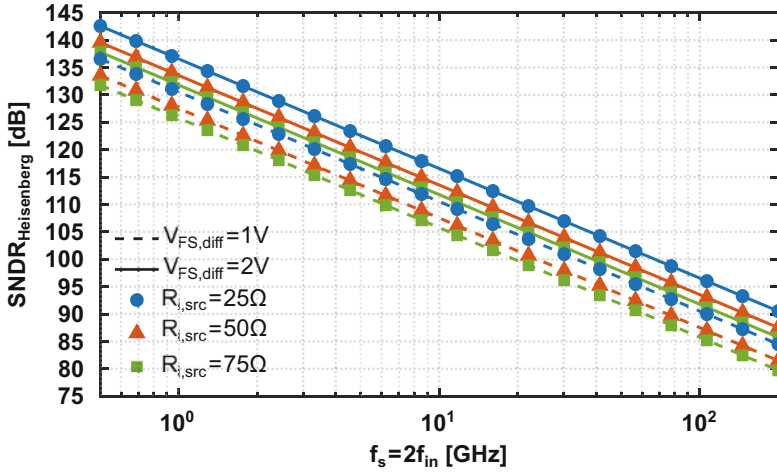


Fig. 2.25 Fundamental accuracy-speed limit due to Heisenberg

The more precisely the energy of a particle in a certain state is known, the greater the uncertainty in the interval of time, in which the particle possesses that particular energy.

The principle is described by the mathematical formula

$$\Delta E \cdot \Delta T \geq \frac{h}{4\pi}, \quad (2.65)$$

where ΔE may be interpreted as the required energy to be within $\pm \text{LSB}/2$ of a quantization level, ΔT is the time required to move from one level to another and assumed half a period of f_s , and $h = 6.62617 \cdot 10^{-34}$ J·s is the Planck constant. Under these assumptions and using $R_{i,\text{src}}$ from the model of Fig. 2.17, the above expression for a differential configuration can be written as

$$\frac{V_{\text{pp-diff}}^2}{2^{2ENOB_{\text{Heis}}} \cdot 8R_{i,\text{src}}} \cdot \frac{1}{(2f_s)^2} \geq \frac{h}{4\pi} \Rightarrow 2^{ENOB_{\text{Heis}}} \cdot f_s \leq \frac{V_{\text{pp-diff}}}{2\sqrt{2hR_{i,\text{src}}}}. \quad (2.66)$$

Finally, the maximum achievable SNDR dictated by the Heisenberg uncertainty principle can be obtained by utilizing Eq. (2.32) (Fig. 2.25)

$$\text{SNDR}_{\text{Heisenberg}} = 6.02 \log_2 \left(\frac{V_{\text{pp-diff}}}{2f_s \sqrt{2hR_{i,\text{src}}}} \right) + 1.76. \quad [\text{dB}] \quad (2.67)$$

2.4.7 Putting It All Together

To finalize our analysis, in Fig. 2.26, we plot all the previously derived accuracy-speed and accuracy-power limits for certain design choices and parameters. As seen in Fig. 2.26, the quantizer metastability for an a_{er} of $1e-5$ is the dominant accuracy limitation when increasing the sample rate above about 25 GS/s. Below this frequency, aperture jitter of 50 fs dominates the accuracy degradation down to about 4 GS/s. At lower sample rates, mismatch is the main limitation to the

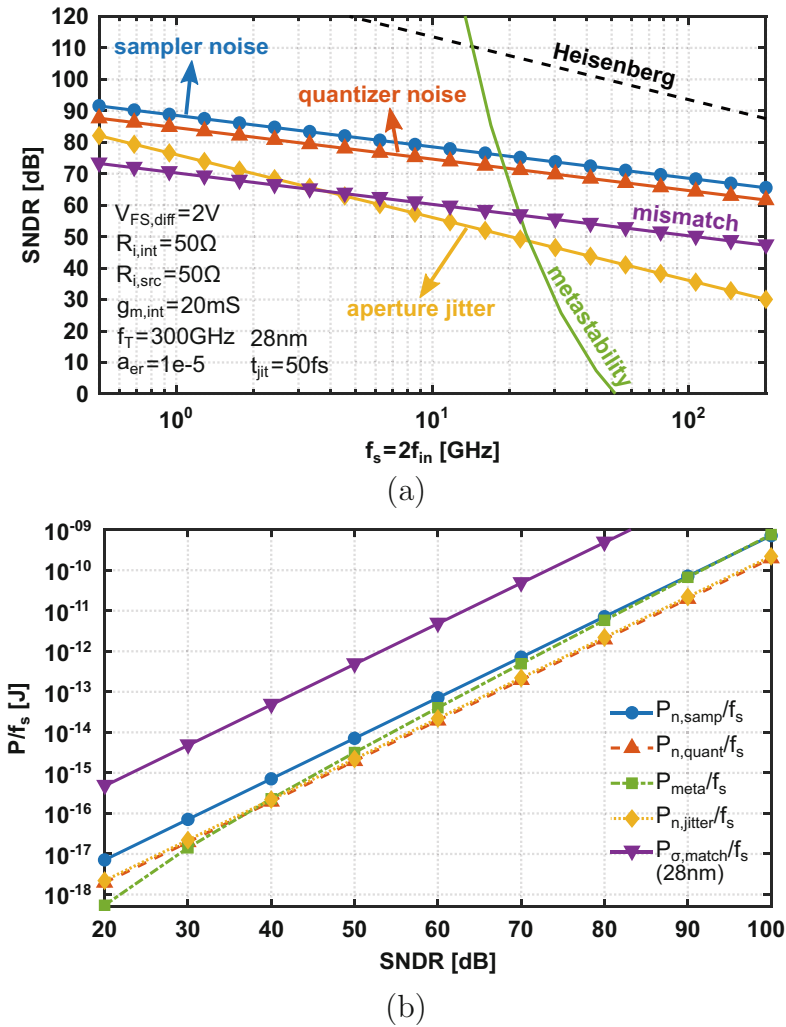


Fig. 2.26 Fundamental limit curves from all the error sources analyzed in this chapter: (a) accuracy-speed and (b) accuracy-power

achievable resolution. Assuming that mismatch is compensated, thermal noise starts limiting the achievable resolution for sample rates below 500 MS/s, with the quantizer as the dominant error source based on our derivations. This is expected at very low sample rates due to the steeper slope of the jitter-limited resolution. The physical Heisenberg uncertainty principle limitation is about 30 dB above the next limitation.

Regarding Fig. 2.26b, our simplified derivations indicate a maximum of about half an order of magnitude power consumption difference between the various noise and metastability limitations. For a 28 nm process, mismatch imposes a power consumption limit of about two orders of magnitude higher than the rest. In reality, the power of the sampler is expected to increase in the presence of an analog front-end with a certain settling requirement. Also, the power estimation for a certain jitter neglects multiple stages in the chain of Fig. 2.22 to realize a certain clock edge steepness, which will inevitably increase this power. Nevertheless, the important take from this first-order power analysis is that every contribution in a converter necessitates an equally careful optimization and/or compensation to yield the best overall results.

2.5 Conclusion

This chapter laid out the fundamental concepts of the A/D conversion process. Its two primary concepts of sampling (time discretization) and quantization (amplitude discretization) were thoroughly discussed. In order to prevent loss of information and yield the sampling process reversible, the Nyquist criterion dictates that the sample rate be at least twice the instantaneous bandwidth of the signal under sampling. The signal may be located in any of the Nyquist zones, and as long as it is band-limited within one, the Nyquist criterion is satisfied. Each sampled value is compared against 2^B discrete levels, and its amplitude is rounded to the nearest level by the quantizer. This rounding process introduces a deterministic quantization error ϵ_q , which under certain conditions can be approximated as white noise, and imposes the ideal single conversion error source. From this analysis, the maximum possible accuracy of a B -bit converter was derived in terms of its SQNR. The major error sources from the circuit blocks in a practical converter chain were identified to deteriorate the performance beyond the quantization error threshold. In the form of noise, these include the sampler thermal noise, the quantizer thermal noise, and the aperture jitter from the clock and input of the sampler. Simple models were introduced, and closed-form expressions were developed to quantify these errors in terms of design parameters. In the form of non-linearity, DNL and INL from the quantizer as well as INL and harmonic distortion from the other blocks in the chain (sampler and a potential front-end) were identified as the main contributors. Generally, any type of non-linearity originates from circuit imperfections and can be minimized either by proper design choices or by calibration, which was briefly overviewed as well. Further, several critical performance evaluation metrics,

including THD, SNR, SNDR, SFDR, as well as the two widely used figures of merit, FoM_W and FoM_S , were briefly discussed.

Equations serving as first-order guidelines were developed, which established the fundamental accuracy-speed-power limits imposed by (1) the sampler noise, (2) the quantizer noise, (3) the quantizer metastability, (4) the aperture jitter, and (5) ultimately physics under certain assumptions. The limits imposed by mismatch were also quantified and compared to the aforementioned ones. The derived equations provided an insight as to what may be ultimately achievable from the elementary building blocks in a converter and what has to be traded-off to maximize the ratio $accuracy \cdot speed \div power$. It was concluded that the contribution from every block needs to be equally carefully optimized and/or compensated to reach the best possible performance. More importantly, this insight allows a better circuit design optimization, avoiding excessive over-design or under-design that could potentially lead to poor power and/or speed performance for a certain accuracy.

Appendix A: Proper FFT Evaluation Setup

Assume we would like to sample at an f_s rate a one-tone sine wave with an input frequency f_{in} and evaluate its frequency spectrum by means of an FFT with N_{FFT} points. The total FFT evaluation time is found as

$$T_{FFT} = N_{FFT} \cdot 1/f_s. \quad (2.68)$$

The resolution bandwidth or FFT bin size is then given by

$$f_{bin} = 1/T_{FFT} = f_s/N_{FFT}. \quad (2.69)$$

For coherent sampling without using windowing, and to avoid spectral leakage, we must ensure an integer number of signal periods N_{PER} . The input frequency is therefore found as

$$f_{in} = N_{PER} \cdot f_{bin} = N_{PER} \cdot f_s/N_{FFT}. \quad (2.70)$$

The same setup can be followed for a two-tone sine wave¹⁵ with input frequencies f_{in1} and f_{in2} as well, provided that these frequencies fall exactly within FFT bins. One of many ways to guaranteeing this is the following:

$$\begin{aligned} f_{in1} &= N_{PER} \cdot f_{bin} - 2 f_{bin} = N_{PER} \cdot f_s/N_{FFT} - 2 f_s/N_{FFT} \\ f_{in2} &= N_{PER} \cdot f_{bin} + 2 f_{bin} = N_{PER} \cdot f_s/N_{FFT} + 2 f_s/N_{FFT}. \end{aligned} \quad (2.71)$$

¹⁵ It can be generalized to an m-tone sine wave with $f_{in1,2,\dots,m}$.

Finally, N_{PER} and N_{FFT} must be relatively prime, meaning that their only positive common divisor is 1. To give a numerical example, for an $f_s = 1$ GS/s and $N_{\text{FFT}} = 1024$, $N_{\text{PER}} = 79$ satisfies the above requirements, leading to an $f_{\text{in}} = 77.1484$ MHz for a one-tone and $f_{\text{in1}} = 75.1953$ MHz and $f_{\text{in2}} = 79.1016$ MHz for a two-tone sine wave, respectively.