



From conceptual spaces to quantum concepts: formalising and learning structured conceptual models

Sean Tull¹ · Razin A. Shaikh^{1,2} · Sara Sabrina Zemljic¹ · Stephen Clark¹

Received: 4 July 2023 / Accepted: 27 October 2023 / Published online: 15 April 2024
© The Author(s), under exclusive licence to Springer Nature Switzerland AG 2024

Abstract

In this article we present a new modelling framework for structured concepts using a category-theoretic generalisation of conceptual spaces, and show how the conceptual representations can be learned automatically from data, using two very different instantiations: one classical and one quantum. A contribution of the work is a thorough category-theoretic formalisation of our framework. We claim that the use of category theory, and in particular the use of string diagrams to describe quantum processes, helps elucidate some of the most important features of our approach. We build upon Gärdenfors' classical framework of *conceptual spaces*, in which cognition is modelled geometrically through the use of convex spaces, which in turn factorise in terms of simpler spaces called *domains*. We show how concepts from the domains of SHAPE, COLOUR, SIZE and POSITION can be learned from images of simple shapes, where concepts are represented as Gaussians in the classical implementation, and quantum effects in the quantum one. In the classical case we develop a new model which is inspired by the β -VAE model of concepts, but is designed to be more closely connected with language, so that the names of concepts form part of the graphical model. In the quantum case, concepts are learned by a hybrid classical-quantum network trained to perform concept classification, where the classical image processing is carried out by a convolutional neural network and the quantum representations are produced by a parameterised quantum circuit. Finally, we consider the question of whether our quantum models of concepts can be considered conceptual spaces in the Gärdenfors sense.

Keywords Conceptual spaces · Category theory · Quantum cognition · Concept learning

1 Introduction

The study of concepts has a long history in a number of related fields, including philosophy, linguistics, psychology and cognitive science (Murphy 2002; Margolis and Laurence 2015). More recently, researchers have begun to consider how mathematical tools from quantum theory can be used to

model cognitive phenomena, including conceptual structure. The general use of quantum formalism in psychology and cognitive science has led to an emerging area called quantum cognition (Aerts 2009; Pothos and Busemeyer 2013). The idea is that some of the features of quantum theory, such as entanglement, can be used to account for psychological data which can be hard to model classically. Examples include ordering effects in how subjects answer questions (Trueblood and Busemeyer 2011) and concept combination (Aerts and Gabora 2005; Tomas and Sylvie 2015).

Another recent development in the study of concepts has been the application of machine learning to the problem of how artificial agents can automatically learn concepts from raw perceptual data (Higgins et al. 2017, 2018). The motivation for endowing an agent with conceptual representations, and learning those representations automatically from the agent's environment, is that this will enable it to reason and act more effectively in that environment, similar to how humans use concepts (Lake et al. 2017). One hope is that the explicit use of concepts will ameliorate some of the negative

✉ Stephen Clark
steve.clark@quantinuum.com

Sean Tull
sean.tull@quantinuum.com

Razin A. Shaikh
razin.shaikh@quantinuum.com

Sara Sabrina Zemljic
sara.zemljic@quantinuum.com

¹ Quantinuum, 17 Beaumont Street, Oxford OX1 2NA, UK

² Department of Computer Science, University of Oxford, Oxford, UK

consequences of the “black-box” nature of neural architectures currently being used in AI.

In this article we present a new modelling framework for concepts based on the mathematical formalism used in quantum theory, and demonstrate how the conceptual representations can be learned automatically from data, using both classical and quantum-inspired models. A contribution of the work is a thorough category-theoretic formalisation of our framework, following Bolt et al. (2019) and Tull (2021). Formalisation of conceptual models is not new (Ganter and Obiedkov 2016), but we claim that the use of category theory (Fong 2019), and in particular the use of string diagrams to describe quantum processes (Coecke and Kissinger 2017), helps elucidate some of the most important features of our approach to concept modelling. This aspect of our work also fits with the recent push to introduce category theory into machine learning and AI more broadly. The motivation is to make deep learning less ad-hoc and less driven by heuristics, by viewing deep learning models through the compositional lens of category theory (Shiebler et al. 2021).

(Murphy 2002, p.1) describes concepts as “the glue that holds our mental world together”. But how should concepts be modelled and represented mathematically? There are many modelling frameworks in the literature, including the *classical theory* (Margolis and Laurence 2022), the *prototype theory* (Rosch 1973), and the *theory theory* (Gopnik and Meltzoff 1997). Here we build upon Gärdenfors’ framework of *conceptual spaces* (Gärdenfors 2004, 2014), in which cognition is modelled geometrically through the use of convex spaces, which in turn factorise in terms of simpler spaces called *domains*.

Our category-theoretic formalisation of conceptual spaces allows flexibility in how the framework is instantiated and then implemented, with the particular instantiation determined by the choice of category. First we show how the framework can be instantiated and implemented classically, by using the formalisation of “fuzzy” conceptual spaces from Tull (2021), and developing a probabilistic model based on Variational Autoencoders (VAEs) (Rezende et al. 2014; Kingma and Welling 2014). Having “fuzzy” probabilistic representations not only extends Gärdenfors’ framework in a useful way, it also provides a natural mechanism for dealing with the vagueness inherent in the human conceptual system, and allows us to draw on the toolkit from machine learning to provide effective learning mechanisms. Our new model—which we call the *Conceptual VAE*—is an extension of the β -VAE from Higgins et al. (2017), with the concepts having explicit labels and represented as multivariate Gaussians in a factored conceptual space.

We use the Spriteworld software (Watters et al. 2019) to generate simple images consisting of coloured shapes of certain sizes in certain positions, meaning our conceptual spaces contain four domains: COLOUR, SIZE, SHAPE and POSITION.

The main question we investigate for the classical model is a representational learning one: can the Conceptual VAE induce factored representations in a latent conceptual space which neatly separates the individual concepts, and under what conditions? Here we demonstrate that, if the system is provided with supervision regarding the domains, and provided with the corresponding four labels for each training instance (e.g. (*blue, small, circle, top*)), then the VAE can learn Gaussians which faithfully represent the colour spectrum, for example. We also show the Conceptual VAE naturally provides a concept classifier, in the form of the encoder, which predicts a Gaussian for an image that can be compared with the induced conceptual representations using the KL divergence.

Our second instantiation of the abstract framework uses a category for describing quantum processes (Coecke and Kissinger 2017). In this case, the images of shapes are represented as *quantum states* in an underlying Hilbert space and concepts are *quantum effects*. Applying a concept effect to an instance state yields a scalar, which we interpret as specifying how well the instance fits the concept. The factoring of the conceptual space is represented naturally in our models through the use of the tensor product as the monoidal product. We choose to implement our quantum model using a hybrid quantum-classical network trained to perform concept classification, where the classical image processing is carried out by a convolutional neural network (Goodfellow et al. 2016, Ch.9) and the quantum representations are produced by a parameterised quantum circuit (Benedetti et al. 2019). Even though the framework is instantiated at a level of abstraction independent of any particular implementation, the use-case we have in mind is one in which the models are (eventually) run on a quantum computer, exploiting the potential advantages such computers may bring. Here the implementation is a classical simulation of a quantum computation.¹

We demonstrate how the training of the hybrid network produces conceptual representations in the Hilbert space which are neatly separated within the domains. We also show how discarding—which produces mixed effects—can be used when the concept to be learned only applies to a subset of the domains, and how entanglement (together with discarding) can be used to capture interesting correlations across domains.

What are some of the main reasons for applying the formalism of quantum theory to the modelling of concepts? First, it provides an alternative, and interesting, mathematical structure to the convex structure of conceptual spaces

¹ Note that we are not making any claims of “quantum supremacy” (Preskill 2012) for the particular set of quantum models that we implement in this article. However, we do anticipate the possibility of quantum models of concepts satisfying our framework which require quantum hardware for their efficient training and deployment, especially as we scale to more realistic datasets and larger quantum circuits.

(see Section 2.7). Second, this structure comes with features which are well-suited to modelling concepts, such as entanglement for capturing correlations, and partial orders for capturing conceptual hierarchies.² Third, the use of the tensor product for combining domains leads to machine learning models with different characteristics to those typically employed in concept learning, such as the Conceptual VAE (i.e. neural networks which use direct sum as the monoidal product plus non-linearities to capture interactions between features) (Havlicek et al. 2019; Schuld and Killoran 2019). The advantages this may bring, especially with the advent of larger, fault-tolerant quantum computers in the future, is still being worked out by the quantum machine learning community, but the possibilities are intriguing at worst and transformational at best.

Note that, in this article, our goal is to set out a novel framework for concept modelling, and demonstrate empirically—with two very different implementations—how concepts can be learned in practice. Further work is required to demonstrate that the framework can be applied fruitfully to data from a psychology lab—which is one of the goals of quantum cognition (Pothos and Busemeyer 2013)—and also to agents acting in (virtual) environments—one of the goals of agent-based AI (Abramson et al. 2020). Note also that no claims are being made here regarding the existence of quantum processes in the brain, only that some cognitive processes can be effectively modelled at an abstract level using the quantum formalism.

The rest of the article is structured as follows. Section 2 provides a thorough category-theoretic formalisation of our modelling framework, using the language of string diagrams to describe the structured models. Section 3 then describes our first instantiation of the framework, which is a novel adaptation of the variational autoencoder. This section also contains experiments showing how Gaussian concept representations can be learned from images of coloured shapes. Section 4 then describes our quantum instantiation, as well as a hybrid implementation applied to the same image data. The hybrid network uses a CNN for the classical image processing and a parameterised quantum circuit for inducing the concept representations (as quantum effects). Finally, Sections 5 and 6 describe related and future work.

2 Formalising conceptual spaces

Gärdenfors' framework of *conceptual spaces* (Gärdenfors 2004, 2014) models conceptual reasoning in both human and artificial cognition. The approach models cognition geometrically, using convex spaces factorised in terms of

² Section 2.6 describes entanglement; we leave the use of partial orders in experiments for future work.

“elementary” spaces called *domains*. Examples include the domains of COLOUR, TASTE, SOUND, and TIME. Concepts are represented as convex regions, or more generally as “fuzzy” functions defined over the space. We begin with a brief formalisation of this framework. While many have been presented (Aisbett and Gibbon 2001; Rickard et al. 2007; Lewis and Lawry 2016; Bechberger and Kühnberger 2017), we draw on the categorical approaches (Bolt et al. 2019; Tull 2021) and the latter's treatment of fuzzy concepts.

Definition 1 A *convex space* is a set Z which forms a measurable space, i.e. is given with a σ -algebra of ‘measurable’ subsets $\Sigma_Z \subseteq \mathbb{P}(Z)$, and which moreover comes with operations which allow us to take convex combinations of elements. That is, for all $p_1, \dots, p_n \in [0, 1]$ with $\sum_{i=1}^n p_i = 1$, we have an operation on Z denoted:

$$(z_1, \dots, z_n) \mapsto p_1 \cdot z_1 + \dots + p_n \cdot z_n$$

These operations are related by axioms one would expect from the familiar example of combinations in a vector space. In particular, the order of elements in a combination doesn't matter, iterated convex combinations are given by multiplying weights, and we always have $1 \cdot z + 0 \cdot z' = z$. For a full formal definition see Bolt et al. (2019).

Definition 2 A *conceptual space* is a convex space Z given as a subset of a product of convex spaces:

$$Z \subseteq Z_1 \times \dots \times Z_n$$

where the product is equipped with element-wise convex operations and the product σ -algebra $\Sigma_{Z_1 \times \dots \times Z_n}$ of the σ -algebras $\Sigma_{Z_1}, \dots, \Sigma_{Z_n}$. We call an element $z = (z_1, \dots, z_n) \in Z$ an *instance* of the conceptual space, following Clark et al. (2021).

Any factor Z_i can be considered a conceptual space itself, with each z_i an instance. A conceptual space is often written as a product of domains, such as COLOUR or SOUND. Each domain itself factorises as a (subset of a) product of *dimensions*. For example, the SOUND domain has the dimensions of PITCH and VOLUME. Here we simply use the neutral term “factor” to treat either dimensions or domains.

Definition 3 A *crisp concept* in a conceptual space Z is a measurable subset $C \subseteq Z$ which is *convex*, meaning it is closed under convex combinations. When $z \in C$ we say z is an *instance* of C .

Convexity means that any point lying “in-between” two instances of a concept will again form an instance of the concept, and is justified by Gärdenfors using experimental evidence in the division of colour space, and the ease of learning convex regions (Gärdenfors 2004).

More generally, it is natural to consider concepts C which are graded or “fuzzy”. To make sense of this, first we will from now on often abuse notation slightly and equate a crisp concept $C \subseteq Z$ with its corresponding indicator function 1_C , writing it as $C: Z \rightarrow [0, 1]$. Then we have $C(z) = 1$ if $z \in C$ and $C(z) = 0$ otherwise.

For a fuzzy concept, the degree $C(z)$ to which z is an instance of a concept should now be able to take any value between 0 (“not at all”) to 1 (“fully satisfied”). In Tull (2021) it is shown that to be well-behaved compositionally and also satisfy a natural generalisation of convexity known as “quasi-concavity”, the membership function should satisfy the following.

Definition 4 A fuzzy concept of Z is a measurable function $C: Z \rightarrow [0, 1]$ which is *log-concave*:

$$C(pz + (1 - p)z') \geq C(z)^p C(z')^{1-p} \tag{1}$$

for all $z, z' \in Z$ and $p \in [0, 1]$. A *prototypical instance* of C is an instance z with $C(z) = \max_{w \in Z} C(w)$, whenever such an instance exists.

The collection of prototypical instances of a fuzzy concept always forms a crisp concept. Conversely, any crisp concept $C \subseteq Z$ may be seen as a special case of a fuzzy concept via its indicator function 1_C as above, with C as its subset of prototypical instances. From now on, for a crisp concept we will not distinguish sharply between the subset and its indicator function, denoting both by C .

Example 1 Any convex subset $Z \subseteq \mathbb{R}^d$ forms a conceptual space, taking Σ_Z to be the Lebesgue measurable subsets. Thus any product $Z = Z_1 \times \dots \times Z_n$ of convex subsets $Z_i \subseteq \mathbb{R}^{d_i}$ forms a conceptual space. In general any product of fuzzy concepts yields a new one on any convex subset $Z \subseteq Z_1 \times \dots \times Z_n$ via:

$$C(z) = \prod_{i=1}^n C_i(z_i) \tag{2}$$

The following example of a fuzzy concept will form the basis of our classical implementation of the framework in Section 3.

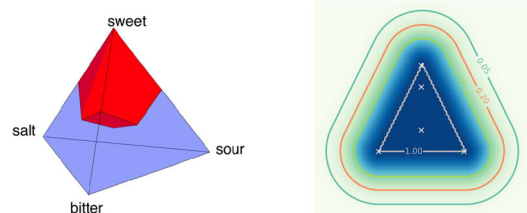
Example 2 We may define a fuzzy concept on $Z = \mathbb{R}^n$ from any multivariate Gaussian with mean μ and covariance matrix Σ :

$$C(z; \mu, \Sigma) = e^{-\frac{1}{2}(z-\mu)^T \Sigma^{-1} (z-\mu)} \tag{3}$$

$$= e^{-\sum_{i=1}^n \frac{1}{2\sigma_i^2} (z_i - \mu_i)^2} \tag{4}$$

In the second line we restrict to the case where Σ is diagonal, with i -th diagonal entry σ_i^2 . In this case C is given as a product of one-dimensional Gaussians $C_i(z_i; \mu_i, \sigma_i^2)$ as in Eq. 2.

Example 3 A simple TASTE domain from Bolt et al. (2019), left-hand below, is given as a convex subset of \mathbb{R}^3 generated by the points *sweet*, *bitter*, *salt* and *sour*. Highlighted in red is a crisp concept for *sweet*. Right-hand below shows a fuzzy concept on \mathbb{R}^2 from Tull (2021). From a set of exemplars (white crosses) the convex closure is formed, yielding the crisp concept P given by the inner triangle. A fuzzy concept is then defined by $C(x) = e^{-\frac{1}{2\sigma^2} d(x, P)^2}$ where $d_H(x, P) = \inf_{p \in P} d(x, p)$, where each point in P is prototypical.

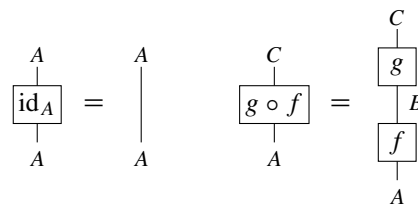


2.1 Categorical setup

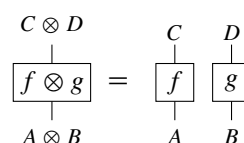
Our aim will now be to lift these basic notions from conceptual space theory into a general categorical framework, allowing us to pass them from the classical to the quantum setting in a principled manner. Here we introduce the categorical preliminaries.

We will work in a *symmetric monoidal category* (\mathbf{C}, \otimes, I) . Recall that this consists of a collection of *objects* A, B, \dots and a collection of *morphisms*, where a morphism from A to B is denoted $f: A \rightarrow B$. We may compose morphisms in sequence when their types match, as well as composing objects and morphisms in parallel via a ‘tensor’ operation \otimes . For more details see e.g. (Coecke 2006).

We will make use of the graphical calculus (Selinger 2010) in which objects are depicted as labelled wires, and morphisms $f: A \rightarrow B$ as boxes with input wire A and output wire B , read here from bottom to top. Identity morphisms and sequential composition are depicted as follows.



Parallel composition via the tensor \otimes is given by drawing diagrams side-by-side.



The (identity on the) monoidal unit I is the empty diagram. Morphisms $\omega: I \rightarrow A, e: A \rightarrow I$ and $r: I \rightarrow I$ are called *states*, *effects* and *scalars* respectively, depicted with no input, output or neither, respectively.

Here we consider categories \mathbf{C} with further structure. First, each object A will come with a distinguished *discarding* effect denoted \ddagger_A , which we interpret as “throwing the system away”, with $\ddagger_I = \text{id}_I$ and $\ddagger_{X \otimes Y} = \ddagger_X \otimes \ddagger_Y$. A morphism $f: A \rightarrow B$ is called a *channel* when it preserves discarding, as in left-hand below. A special case is that we call a state ω of A *normalised* when the right-hand equation below holds.



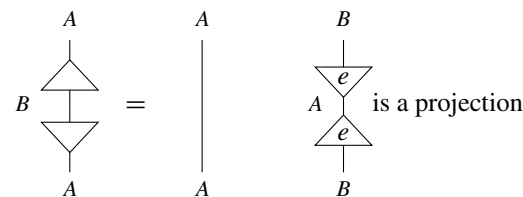
We also assume \mathbf{C} is enriched in partial orders, so that each homset $\mathbf{C}(A, B)$ forms a partially ordered set \leq , so that we may now have $f \leq g$ for a pair of morphisms $f, g: A \rightarrow B$. Moreover this ordering is ‘respected’ by composition so that if $f \leq g$ then any composite of both sides with the same morphism h still satisfies this relation³.

This allows us to generalise inclusions of convex subsets via the following, related to “comprehensions” (Cho et al. 2015) and “compression” maps in quantum reconstructions (Tull 2019, Chap. 4).

Definition 5 A *projection* is a morphism $p: A \rightarrow A$ with $\ddagger_A \circ p \leq \ddagger_A$ and such that:

1. For all $f: A \rightarrow B$ we have $\ddagger_B \circ f \leq \ddagger_A \circ p \implies f = f \circ p$;
2. For all $f: B \rightarrow A$ we have $\ddagger_A \circ f \leq \ddagger_A \circ p \circ f \implies f = p \circ f$.

It follows that $p = p \circ p$. An *embedding* of an object A into B is given by a channel $e: A \rightarrow B$ and morphism $e^\dagger: B \rightarrow A$, depicted using triangles as below, such that $e^\dagger \circ e = \text{id}_A$ and $p = e \circ e^\dagger$ is a projection.



We often simply call the morphism e the *embedding*, and e^\dagger the *projection*, of the pair.

³ For example if $f \leq g$ then $h \circ f \leq h \circ g, f \circ h \leq g \circ h$ and $f \otimes h \leq g \otimes h$ where h is any morphism h of an appropriate type for each case.

Any channel which is an isomorphism $A \simeq B$ forms a special case of an embedding, where $e^\dagger = e^{-1}$. Another important special case is an embedding of I into A , which we call a *point* of A .⁴ By definition it includes a normalised state ψ with an effect $\psi^\dagger \leq \ddagger$ satisfying:

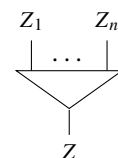
Embeddings are always closed under composition in the following sense.

Lemma 6 If $d: A \rightarrow B$ and $e: B \rightarrow C$ are embeddings then so is $e \circ d: A \rightarrow C$, with projection $d^\dagger \circ e^\dagger$.

2.2 Conceptual models

Let us now see how each of our earlier features from conceptual space theory can be described in a general category \mathbf{C} with the structure outlined in Section 2.1. Firstly, monoidal categories immediately allow us to describe the compositions of factors Z_i appearing in a conceptual space, as follows.

Definition 7 A *conceptual model*⁵ is given by an object Z along with an indexed collection of objects Z_1, \dots, Z_n , called the *factors*, and an embedding of Z into $Z_1 \otimes \dots \otimes Z_n$.



For simplicity we refer to a model as Z , with the factors and embedding implicit. Often the embedding is an isomorphism $Z \simeq Z_1 \otimes \dots \otimes Z_n$ exhibiting Z as a product of the factors.

Definition 8 A *concept* of a conceptual model Z is an effect C on Z .



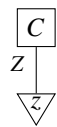
⁴ Later we will define instances as special cases of points. Instances and points differ in quantum models, because of entanglement, but coincide classically.

⁵ Henceforth we use the generic term “model” rather than “space” since a conceptual model can be defined in a category without any spatial character.

An *instance* is a point z of Z which forms a product of points z_i of the factors Z_i , as below:

$$(6)$$

The order structure on morphisms means that the concepts are automatically partially ordered. We interpret $C \leq D$ as stating that D is a “more general” concept than C . The factorisation property Eq. 6 generalises the fact that in a conceptual space every instance $z = (z_1, \dots, z_n)$ factors as a product of one instance z_i per factor Z_i . Composing a concept C with any input state, in particular any instance z , will yield a scalar. For an instance we interpret this as specifying how well the instance fits the concept:



We say that an instance z is *prototypical* for a concept C when $C \circ w \leq C \circ z$ for all instances w . It remains for us to identify those concepts which are crisp.

Definition 9 A concept C on Z is *crisp* when it is of the form

$$C = K$$

for some projection morphism $Z \rightarrow K$ induced by an embedding of K into Z . When the projection is given by a point of Z we call C a *pure* concept.

By definition each crisp concept has $C \leq \ddagger$. Intuitively we can identify the crisp concept with object K via its embedding e . Indeed by the definition of an embedding, for any instance z of Z we have $C \circ z = 1$ iff $z = e \circ k$ for some point k of K . Moreover any concept D with $D \leq C$ restricts to K in that $D = E \circ e^\dagger$ for some effect E on K . A pure concept can be thought of as a “maximally sharp” concept, being of the form $z = \psi^\dagger$ as in Eq. 5 where $z = \psi$ is in fact a point of Z .

2.3 Classical conceptual models

Let us now meet our main classical examples of categories and their notions of conceptual model.

Class: discrete probability In the category **Class** the objects are finite sets and the morphisms $M: X \rightarrow Y$ are functions $M: Y \times X \rightarrow \mathbb{R}^+$. We think of such a morphism

as a ‘matrix’ with values in \mathbb{R}^+ , indexed by the elements of Y, X , which we write as $(M(y | x))_{y \in Y, x \in X}$. Composition is matrix multiplication:

$$(N \circ M)(z | x) := \sum_{y \in Y} N(z | y)M(y | x)$$

Identity morphisms satisfy $\text{id}_X(y | x) = \delta_{x,y}$. $X \otimes Y = X \times Y$, with $I = \{\star\}$ the singleton set, and $M \otimes N$ the Kronecker product of matrices. We can equate states ω and effects e of X each with functions $X \rightarrow \mathbb{R}^+$. In particular, scalars correspond to positive reals $s \in \mathbb{R}^+$. \ddagger is the function $x \mapsto 1$ for all $x \in X$.

A state ω of X is normalised iff it describes a probability distribution, with $\sum_{x \in X} \omega(x) = 1$. More generally, a morphism $M: X \rightarrow Y$ is a channel iff it is a finite probability channel (*Stochastic matrix*) with $\sum_{y \in Y} M(y | x) = 1$ for each $x \in X$. \leq is the element-wise ordering from \mathbb{R}^+ . The points of X are precisely the point distributions δ_x for $x \in X$. An embedding $X \hookrightarrow Y$ is given by an inclusion of a subset $X \subseteq Y$ via $x \mapsto \delta_x$, and its projection $Y \rightarrow X$ is given by $y \mapsto \delta_y$ when $y \in X$ and $y \mapsto 0$ otherwise.

A conceptual model in **Class** is thus a finite set Z given as a subset $Z \subseteq Z_1 \times \dots \times Z_n$. A concept is an arbitrary function $C: Z \rightarrow \mathbb{R}^+$, ordered point-wise. An instance is any element $z = (z_1, \dots, z_n) \in Z$, with Eq. 6 holding automatically. Applying a concept C to an instance z gives $C(z) \in \mathbb{R}^+$. Crisp concepts are the indicator functions 1_K of arbitrary subsets $K \subseteq Z$, while pure concepts are those of instances $z \in Z$.

Prob: measure-theoretic probability In **Prob** the objects are measurable spaces (X, Σ_X) . A morphism $f: X \rightarrow Y$ is a *Markov (sub)kernel*, a function sending each $x \in X$ to a sub-probability measure $f(x)$ over Y , in a “measurable” way (Panangaden 1998; Cho and Jacobs 2019). Composition of $f: X \rightarrow Y$ and $g: Y \rightarrow Z$ is via integration⁶:

$$(g \circ f)(x, A) := \int_{y \in Y} g(y, A)df(x)(y)$$

for each $x \in X, A \in \Sigma_Z$. The identity sends each x to the point measure δ_x . We set $X \otimes Y = X \times Y$, with I being the singleton set, and define $f \otimes g$ to send each pair (x, y) to the *product measure* of the measures $f(x)$ and $g(y)$. States of X may be identified with sub-probability measures ω over X , and are normalised iff they form a probability measure, with $\omega(X) = 1$. Effects correspond to measurable functions $e: X \rightarrow [0, 1]$. \ddagger is the constant function at 1. Scalars are probabilities $p \in [0, 1]$. Composing a state with an effect yields the expectation value $e \circ \omega = \int_{x \in X} e(x)d\omega(x) \in \mathbb{R}^+$.

⁶ Here we use the standard definition of integration on a measurable space, which exists since $g(-, A)$ is measurable and bounded in $[0, 1]$ by assumption.

A morphism $f: X \rightarrow Y$ is a channel iff it sends each $x \in X$ to a probability measure. Then $f \leq g$ whenever $f(x, A) \leq g(x, A)$ for all $x \in X, A \in \Sigma_Y$. An embedding $X \hookrightarrow Y$ is an inclusion of a subset $X \subseteq Y$ via $x \mapsto \delta_x$ for $x \in X$, with the projection $Y \rightarrow X$ given by $y \mapsto \delta_y$ when $y \in X$ and $y \mapsto 0$ otherwise.

A conceptual model in **Prob** is thus a measurable space given as a measurable subset $Z \subseteq Z_1 \times \dots \times Z_n$ of spaces Z_i . Concepts are measurable functions $C: Z \rightarrow [0, 1]$, instances and pure concepts correspond to points $z \in Z$, crisp concepts 1_K correspond to arbitrary measurable subsets $K \subseteq Z$.

ConSp: conceptual spaces The category **ConSp** (Tull 2021) is defined just like **Prob** except that the objects are now convex spaces and morphisms are (sub)kernels f which are *log-concave*, meaning that

$$f(px + (1 - p)y, pA + (1 - p)B) \geq f(x, A)^p f(y, B)^{1-p} \tag{7}$$

for all $p \in [0, 1], x, y \in X$ and $A, B \in \Sigma_Y$. Here $X \otimes Y = X \times Y$ is the product of convex spaces, with element-wise convex operations.

A conceptual model in **ConSp** is precisely a conceptual space, i.e. a convex space viewed as a convex subset $Z \subseteq Z_1 \times \dots \times Z_n$ of convex spaces Z_i . Instances are points $z \in Z$. Crisp concepts are precisely those of Definition 3, namely the indicator functions 1_K of convex measurable subsets $K \subseteq Z$, with pure concepts being the indicator functions 1_z of points $z \in Z$. Concepts are fuzzy concepts $C: Z \rightarrow [0, 1]$ in the sense of Definition 4.

2.4 Quantum conceptual models

We can now define our quantum model of concepts inspired by the conceptual space framework. To do so we will simply unpack our definitions from Section 2.2 in the following category of quantum processes.

Quant: Quantum processes In the category **Quant** the objects are finite dimensional Hilbert spaces, and morphisms $f: \mathcal{H} \rightarrow \mathcal{K}$ are *completely positive* (CP) maps $f: L(\mathcal{H}) \rightarrow L(\mathcal{K})$, where $L(\mathcal{H})$ is the space of linear operators on \mathcal{H} . Such a map f is linear, and such that for any \mathcal{H}' the map $g = f \otimes \text{id}_{\mathcal{H}'}$ is *positive* in that whenever a is a positive operator then $g(a)$ is also. We set $f \leq g$ whenever $g - f$ is CP.

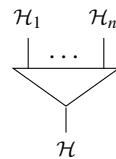
Here \otimes is the usual tensor of Hilbert spaces and linear maps, with $I = \mathbb{C}$. In particular, states ω and effects e on \mathcal{H} may both be identified with positive operators $a \in L(\mathcal{H})$ via $a = \omega(1)$ and $e(b) = \text{Tr}(ab)$, respectively, where Tr denotes the trace. Scalars are again positive reals $r \in \mathbb{R}^+$. Discarding is the functional $\sharp(a) = \text{Tr}(a)$, corresponding to the identity operator $\text{id}_{\mathcal{H}}$.

A morphism f is a channel iff it is a completely positive trace-preserving (CPTP) map, with $\text{Tr}(f(a)) = \text{Tr}(a)$ for all $a \in L(\mathcal{H})$. A state ρ is normalised precisely when it is a density matrix, with $\text{Tr}(\rho) = 1$.

A special class of morphisms are the *pure* CP maps $\hat{f}: L(\mathcal{H}) \rightarrow L(\mathcal{K})$, given by $\hat{f}(a) = f \circ a \circ f^\dagger$ for some linear map $f: \mathcal{H} \rightarrow \mathcal{K}$. All other morphisms are called *mixed*. Embedding morphisms are the pure maps induced by inclusions $i: \mathcal{K} \hookrightarrow \mathcal{H}$ of subspaces into \mathcal{H} . The corresponding projection is the pure map induced by the linear projection i^\dagger onto \mathcal{K} . A point of \mathcal{H} may be identified with a pure quantum state $|\psi\rangle\langle\psi|$ for some unit vector $\psi \in \mathcal{H}$.⁷

We now arrive at our quantum adaptation of the conceptual space framework.

Definition 10 A *quantum conceptual model* is a conceptual model in **Quant**:



Thus a quantum conceptual model is a Hilbert space \mathcal{H} given as a subspace of a tensor product of Hilbert spaces $\mathcal{H} \subseteq \mathcal{H}_1 \otimes \dots \otimes \mathcal{H}_n$. A quantum concept is then precisely a quantum effect, i.e. a positive operator $C \in L(\mathcal{H})$, ordered via $C \leq D$ whenever $D - C$ is positive. An instance is a pure state $|\psi\rangle\langle\psi|$ given by a unit vector $\psi \in \mathcal{H}$, which furthermore factorises as

$$\psi = \psi_1 \otimes \dots \otimes \psi_n \tag{8}$$

for unit vectors $\psi_i \in \mathcal{H}_i$, giving it a well-defined pure state value on each factor \mathcal{H}_i . All instances are pure, with mixed states ρ interpreted as states of uncertainty (i.e. probabilistic mixtures) over pure states such as instances. In contrast concepts may be mixed or pure. The application of a quantum concept C to an instance ψ is given by

$$\begin{array}{c} \boxed{C} \\ \downarrow \\ \mathcal{H} \end{array} = \langle \psi | C | \psi \rangle \in \mathbb{R}^+$$

More generally applying C to a mixed state ρ yields $\text{Tr}(C\rho) \in \mathbb{R}^+$.

Crisp concepts correspond to subspaces $\mathcal{K} \subseteq \mathcal{H}$. More precisely, any such subspace defines a crisp concept via the

⁷ Here we use the standard “bra-ket” notation whereby vectors and linear functionals on \mathcal{H} are written in the form $|\psi\rangle, \langle\phi|$ respectively. Then for a unit vector $\psi \in \mathcal{H}$, $|\psi\rangle\langle\psi|$ is the density operator of the corresponding pure state on \mathcal{H} .

projection operator P onto \mathcal{K} with $P(\psi) = \psi$ for ψ in \mathcal{K} and $P(\psi) = 0$ for ψ in \mathcal{K}^\perp .

Pure quantum concepts are precisely those crisp quantum concepts which are themselves pure as effects. For these, \mathcal{K} is given by a one-dimensional subspace $\langle \psi \rangle$ spanned by some unit vector $\psi \in \mathcal{H}$. Thus pure quantum concepts are precisely effects of the form $|\psi\rangle\langle\psi|$ where ψ is any unit vector (not necessarily an instance). Such a concept sends each instance ϕ to $|\langle\psi|\phi\rangle|^2 \in [0, 1]$.

Example 4 A quantum conceptual model $\mathcal{H} = \mathcal{H}_{\text{Hue}} \otimes \mathcal{H}_{\text{Sat}} \otimes \mathcal{H}_{\text{Light}}$ for COLOUR with factors HUE, SATURATION and LIGHTNESS is given in Yan et al. (2021), where HUE is encoded on a single qubit, represented on the Bloch sphere. In particular each instance (colour) is taken to be a tensor of pure states over each of the factors.

We will meet further examples of quantum conceptual models in Section 4.

2.5 Entangled concepts

It is natural to wonder what advantages, if any, quantum concepts might possess over classical ones. One feature distinguishing quantum models from classical ones is the presence of pure *entangled* concepts. For the following we restrict to categories with scalars given by \mathbb{R}^+ , as in all of our examples here.

Definition 11 A concept C is a *product* concept when there are effects C_1, \dots, C_n such that

$$\begin{array}{c} \boxed{C} \\ | \\ Z \end{array} = \begin{array}{c} \boxed{C_1} \quad \boxed{C_n} \\ | \quad \dots \quad | \\ \text{---} \text{---} \text{---} \\ | \\ Z \end{array} \tag{9}$$

A concept C is *separable* when its value on instances is equal to a convex mixture of product concepts. That is, there are product concepts $C^{(1)}, \dots, C^{(k)}$ such that $C \circ z = \sum_{j=1}^k C^{(j)} \circ z$ for all instances z , where the sum is taken in \mathbb{R}^+ . If a concept C is not separable we say that it is *entangled*.

A product concept treats the factors independently, applying a fixed concept to each. Entangled concepts capture correlations between factors which cannot be reduced to any mixture over such product concepts. **Class**, **Prob** and **ConSp** contain product concepts as well as separable (but non-product) concepts. Nonetheless in **Class** every concept is separable. However, these categories do not contain any pure entangled concepts, since every point of a model $Z \subseteq Z_1 \times \dots \times Z_n$ forms an instance $z = (z_1, \dots, z_n)$ and

hence every pure concept is a product of pure effects z_i^\dagger on each factor.

In contrast, quantum models \mathcal{H} contain both entangled and pure entangled concepts. For any unit vector $\psi \in \mathcal{H}$ which is entangled in the usual sense, i.e. not of the form Eq. 8, the point $|\psi\rangle\langle\psi|$ is not an instance, and its corresponding pure concept on \mathcal{H} is entangled.

Example 5 Consider a Hilbert space \mathcal{H} with orthonormal basis $\{|i\rangle\}_{i=0}^{n-1}$. An entangled pure concept on $\mathcal{H} \otimes \mathcal{H}$ is given by the *Bell effect*, induced by the (unnormalised) vector $\sum_{i=0}^{n-1} |i\rangle |i\rangle$ (where sum denotes superposition), with operator $\sum_{i,j=0}^{n-1} |i\rangle\langle i| \langle j| \langle j|$.

Remark 1 The finite sum in Definition 11 should ultimately be replaced with an integral, so that each concept in **Prob** is separable. It would be interesting to explore whether entanglement exists in **ConSp**.

2.6 Quantum and classical concept combinations

To compare classical and quantum concepts, and to demonstrate the role of entangled concepts in quantum models, let us now consider the ways in which we may “combine” (crisp) concepts in each of our example categories. Given a collection of crisp concepts $(C_i)_{i=1}^n$, by a *combination* we mean a new (crisp) concept C such that every prototypical instance of one of the C_i is a prototypical instance of C .⁸

We will focus in particular on the natural scenario in which we are given a model Z and wish to combine (the pure concepts induced by) a collection of instances z_1, \dots, z_n . The result is a concept C with the z_1, \dots, z_n as prototypical instances, which we think of as learned from these exemplars.

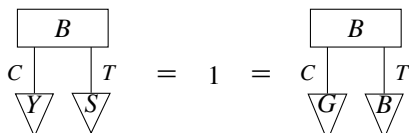
A starting point is to observe that crisp concepts in each category are closed under intersections $\bigcap_{i \in I} C_i$ (of arbitrary, measurable, convex, linear subsets respectively). They hence form a complete lattice with top element \ddagger (and so may be viewed as a *Formal Concept Lattice* in the sense of Ganter and Wille (1999)). This means that one way to combine crisp concepts is via their *disjunction* or least upper bound $C = \bigvee_{i \in I} C_i$.

Classical combinations In **Class** and **Prob**, the disjunction is given by the union of subsets C_i . In fact this is seemingly the only natural way to combine concepts. Indeed here any crisp concept may be identified with its set of prototypical instances, so that any combination C satisfies $\bigvee_{i=1}^n C_i \leq C$. In particular the classical combination $z_1^\dagger \vee \dots \vee z_n^\dagger$ of instances z_1, \dots, z_n is the subset $\{z_1, \dots, z_n\}$.

⁸ In this article “combination” of concepts is always meant in this sense. However there are many distinct meaningful operations on concepts which could also be called their combination, such as the more conjunction-like notion of combining “pet” and “fish” into “pet fish” (Aerts and Gabora 2005).

Spatial combinations In **ConSp** the disjunction is given by the *convex closure* $C = \text{Conv}(\bigcup_{i \in I} C_i)$ of the convex subsets C_i , the smallest convex subset containing all of them. Again any combination C has $\bigvee_{i=1}^n C_i \leq C$. The spatial combination of z_1, \dots, z_n now includes any convex combination of them.

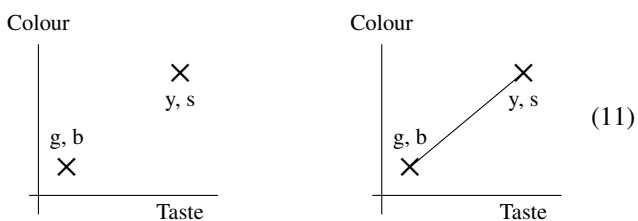
Example 6 Consider a model with factors $C = \text{COLOUR}$ and $T = \text{TASTE}$ and a concept B for *banana* which combines two instances: a yellow (Y) sweet (S) banana, and a green (G) bitter (B) banana.



For simplicity, suppose yellow and green are “orthogonal” in that $Y^\dagger \circ G = 0$. The classical combination yields the crisp concept whose only points are the two instances $(Y, S), (G, B)$ themselves, which by orthogonality can be equivalently written as a sum using element-wise addition of matrices in **Class**:

$$\begin{array}{c} \boxed{D} \\ \downarrow \quad \downarrow \\ C \quad T \end{array} = \begin{array}{c} \triangle Y \\ \downarrow \\ C \end{array} \begin{array}{c} \triangle S \\ \downarrow \\ T \end{array} + \begin{array}{c} \triangle G \\ \downarrow \\ C \end{array} \begin{array}{c} \triangle B \\ \downarrow \\ T \end{array} \tag{10}$$

The classical combination is depicted left-hand below. The spatial combination instead corresponds to the line connecting the two points (right-hand below).



Quantum combinations In **Quant**, the disjunction is given by the linear closure $C = \text{Lin}(\bigcup_{i \in I} C_i)$ of all the subspaces C_i , the smallest subspace containing all of them. This yields a mixed quantum concept which we may interpret as their “coarse-graining”, and again refer to as their *classical combination*. Crucially, however, in a quantum conceptual model there are in fact *many* possible ways to combine crisp concepts, aside from the disjunction, even into a pure concept. That is, there are combinations C of the crisp concepts C_i which do not satisfy $\bigvee_i C_i \leq C$.

Definition 12 In a quantum conceptual model, by a *quantum combination* of instances ψ_1, \dots, ψ_n we mean a pure concept $C = \phi^\dagger$ with these instances as prototypical.

The presence of quantum combinations is closely related to entanglement, coming from the fact that instances are only a subset of the points in a quantum model, since they are non-entangled. Indeed any quantum combination of two or more instances will be entangled.

Example 7 The Bell effect in Example 5 is a pure concept with prototypical instances being precisely those of the form $|\psi^*\rangle \otimes |\psi\rangle$ for unit vectors $|\psi\rangle$, where $|\psi^*\rangle$ denotes the conjugate vector with respect to the given basis. Thus it forms a pure quantum combination of any such instances.

Example 8 Consider again the setting of the *banana* concept combination from Example 6. In **Quant** we can form the classical combination of instances which is again of the form Eq. 10, where $+$ is now the sum of CP maps. Alternatively, we may form a quantum combination $|\psi\rangle\langle\psi|$ where:

$$\psi = |Y, S\rangle + |G, B\rangle \in C \otimes T \tag{12}$$

More generally any linear map $f: C \rightarrow T$ such that $f(|Y\rangle) = |S\rangle, f(|G\rangle) = |B\rangle$ defines a suitable entangled concept $E = \frown \circ (f \otimes \text{id})$, where \frown denotes the Bell Effect from Example 5. Consider the case where $C = T = \mathbb{C}^2, |Y\rangle = |S\rangle = |0\rangle$ and $|G\rangle = |B\rangle = |1\rangle$. A quantum combination E is now given by the Bell effect. The classical and quantum combinations D, E act on instances as follows:

$$\begin{array}{c} \boxed{D} \\ \downarrow \quad \downarrow \\ \triangle \psi \quad \triangle \phi \end{array} = \sum_{i=0}^1 |\langle i|\psi\rangle|^2 |\langle i|\phi\rangle|^2$$

$$\begin{array}{c} \boxed{E} \\ \downarrow \quad \downarrow \\ \triangle \psi \quad \triangle \phi \end{array} = |\langle \psi^* | \phi \rangle|^2$$

The classical combination D simply compares any input to the two instances, with no further prototypical instances besides those given. As a result the structure of each space “between” $|0\rangle$ and $|1\rangle$ is lost, with the orthogonal states $|\pm\rangle = \frac{1}{\sqrt{2}}(|0\rangle \pm |1\rangle)$ treated identically and $D(|+\rangle \otimes |-\rangle) = \frac{1}{2}$. In contrast the quantum combination E can be seen to encode a structural relationship between the factors, generalising from $|00\rangle, |11\rangle$. Any instance $|\phi^*\rangle \otimes |\phi\rangle$ is prototypical, and conversely, tensors of (conjugate) orthogonal points will not fit the concept, e.g. $E(|+\rangle \otimes |-\rangle) = 0$.

In the above example we see that entangled quantum concept combinations can encode relationships between factors, rather than simply (weighted) collections of exemplars. Indeed any pure entangled concept on $C \otimes T$ corresponds to a pure linear map $f: C \rightarrow T$. We can understand this as a generalisation from the instances into a structural relationship between the factors, akin to a concept of the form

$\{(x, f(x)) \mid x \in C\}$ where f is now affine (convexity-preserving).

As such, quantum combinations share the benefits of spatial combinations on a conceptual space, in that one may form structured concepts by generalising from a small set of instances, as on the right of Eq. 11. However, in the quantum case this can be encoded even within a *single* pure concept. Our conclusion is that entanglement provides an effective way for concepts to encode relationships between factors in quantum conceptual models.

2.7 Is a quantum model a conceptual space?

In comparing conceptual spaces with quantum models, it is natural to ask whether we may view the latter as an instance of the former, while our discussion of entangled concepts in the previous section suggested they should be considered distinct. We now discuss this question in detail. We begin with the case of a model with only a single factor, described by a Hilbert space \mathcal{H} .

Hilbert space as a convex space Naively we can first observe that, as a complex vector space, \mathcal{H} does indeed count as a convex space according to Definition 1. However, arbitrary vectors in \mathcal{H} do not have a direct physical interpretation as states, but only the unit vectors ψ (after identification up to global phase) which form the pure states. These pure states do not straightforwardly form a convex space in the sense of Definition 1, since convex combinations of unit vectors are not unit vectors in general.

Pure states as a betweenness space We can nonetheless view the pure states as a geometric space, in a different way. This is most evident for a qubit $\mathcal{H} = \mathbb{C}^2$, whose pure states are visualised via the surface of the *Bloch sphere*. Though the surface of the sphere does not come with convex mixing in the sense of Definition 1, it forms an instance of a broader notion of convex space which may be used to formalise conceptual spaces, known as a *Betweenness space* (Gärdenfors 2004, 2014; Aisbett and Gibbon 2001). This is a set Z along with a ternary operation $B(x, y, z)$ which says that the point y is “in-between” x and z . A subset S is then *convex* if whenever $x, z \in S$ and $B(x, y, z)$, then $y \in S$ also. The Bloch sphere forms a Betweenness space when defining $B(x, y, z)$ whenever a geodesic from x to z passes through y ; see Fig. 1.

We now ask: is the quantum model of concepts on \mathbb{C}^2 the same as that given by the Bloch sphere as a Betweenness space? In fact the sets of concepts in each model are distinct. Firstly, crisp concepts in the quantum model correspond to subspaces, which on the Bloch sphere are either single points (dimension 1) or the entire sphere (dimension 2). So most convex regions on the sphere, the crisp concepts in the Betweenness space Z , are not valid quantum concepts. Conversely, most quantum concepts are not valid

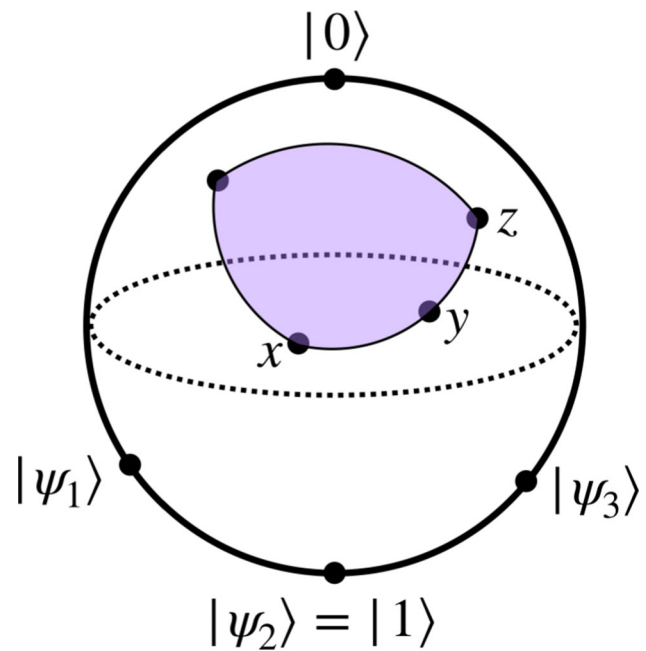


Fig. 1 The Bloch sphere as a Betweenness space, with marked examples of betweenness $B(x, y, z)$, and a convex region shown in purple. The states $|\psi_i\rangle\langle\psi_i|$ are used to show that $|0\rangle\langle 0|$ is not quasi-concave

fuzzy concepts in the Betweenness space Z . As argued in Tull (2021) and mentioned before Definition 4, a fuzzy concept $C: Z \rightarrow [0, 1]$ should at least satisfy the notion of *quasi-concavity*, which states that if $C(x), C(z) \geq t$ then the same holds for any y with $B(x, y, z)$. Example 9 below demonstrates that quantum concepts can fail to satisfy this condition.

Example 9 Consider the pure concept $C = |0\rangle\langle 0|$. Let $|\psi_i\rangle = \cos(\frac{\theta_i}{2})|0\rangle + \sin(\frac{\theta_i}{2})|1\rangle$ for $i = 1, 2, 3$, as in Fig. 1. Setting $\theta_1 = \frac{2\pi}{3}, \theta_2 = \pi, \theta_3 = \frac{4\pi}{3}$ then $|\psi_2\rangle\langle\psi_2| = |1\rangle\langle 1|$ is between $|\psi_1\rangle\langle\psi_1|$ and $|\psi_3\rangle\langle\psi_3|$, making C not quasi-concave, since

$$C(|\psi_1\rangle\langle\psi_1|) = C(|\psi_3\rangle\langle\psi_3|) = \frac{1}{4} > 0 = C(|\psi_2\rangle\langle\psi_2|).$$

Spaces of mixed states One may be tempted to instead view a quantum conceptual model as a different convex space, namely the space $Z = \text{St}(\mathcal{H})$ of (pure and mixed) density matrices on \mathcal{H} , so that these form the instances $z \in Z$. Indeed it follows from linearity that quantum concepts C do satisfy quasi-concavity on this space. However, since density matrices are interpreted as states of uncertainty over pure quantum states, it is more natural to view them as the analogues of *distributions over* a conceptual space, rather than instances. Finally, even if one attempts to view a quantum model as a convex space $\text{St}(\mathcal{H})$, the manner in which we compose such models via the tensor is fundamentally different, making both classes of models distinct:

$$\text{St}(\mathcal{H} \otimes \mathcal{K}) = \text{St}(\mathcal{H}) \otimes \text{St}(\mathcal{K}) \neq \text{St}(\mathcal{H}) \times \text{St}(\mathcal{K})$$

In summary, due to their treatment of entangled concepts and the arguments above, it is most natural to view quantum models as distinct from conceptual spaces. Nonetheless they possess the same benefits for learnability, replacing convex by linear subspaces, and thanks to entanglement may be even more natural for describing correlated concepts.

3 Classical implementation: the conceptual VAE

Our first implementation comes from instantiating the framework using the **ConSp** category from Section 2.3, and implementing fuzzy concepts as Gaussians, as described in Example 2. There is already an existing literature on learning Gaussian representations of concepts, using a tool from machine learning called the *Variational Autoencoder (VAE)* (Higgins et al. 2017). Here we show how to extend that work by defining a new VAE model which provides explicit representations of concepts which fit our framework.

3.1 VAEs for concept modelling

The Variational Autoencoder (VAE) (Kingma and Welling 2014; Rezende et al. 2014) provides a framework for the generative modeling of data, where the data potentially lives in some high-dimensional space. It uses the power of neural networks to act as arbitrary function approximators to capture complex dependencies in the data (e.g. between the pixels in an image). The VAE uses a latent space \mathbf{Z} which acts as a bottleneck, compressing the high-dimensional data into a lower dimensional space.⁹ The question we investigate is whether the VAE model can be adapted so that \mathbf{Z} has desirable properties from a conceptual space perspective, such as interpretable dimensions which contain neatly separated, labelled concepts from individual domains. First we describe the standard VAE model before describing how to adapt it in order to incorporate labelled concepts.

3.1.1 The vanilla VAE

Figure 2 (left) shows the graphical model for the VAE. In terms of the generative story, which is represented by the solid arrows in the plate diagram, first a point \mathbf{z} in the latent space \mathbf{Z} is sampled according to the prior $p(\mathbf{z})$, and then a data point \mathbf{x} is generated according to the likelihood $p_{\theta}(\mathbf{x}|\mathbf{z})$. The dashed arrows denote the variational approximation $q_{\phi}(\mathbf{z}|\mathbf{x})$ to the intractable posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$. The prior is assumed to be a centered isotropic multivariate Gaussian

$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{1})$ (Kingma and Welling 2014). The approximate posterior $q_{\phi}(\mathbf{z}|\mathbf{x})$ is also assumed to be a multivariate Gaussian with a diagonal covariance matrix, but with means and variances predicted by a neural network with learnable parameters ϕ . In our case, since \mathbf{X} is a dataset of images, q_{ϕ} will be instantiated by a convolutional neural network (CNN), which is referred to as the *encoder*. Similarly, p_{θ} will be instantiated by a de-convolutional neural network (de-CNN), and referred to as the *decoder*.

The function that is optimised during training is the RHS of the following equation (Doersch 2016):

$$\begin{aligned} \log p(\mathbf{x}) - \mathcal{D}(q(\mathbf{z}|\mathbf{x}), p(\mathbf{z}|\mathbf{x})) \\ = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}(q(\mathbf{z}|\mathbf{x}), p(\mathbf{z})) \end{aligned} \quad (13)$$

where \mathcal{D} is the KL divergence. Note that, since the KL on the LHS is positive, the equation provides a lower bound on the likelihood, known as the *evidence lower bound (ELBO)*. The advantage of this formulation is that the RHS can be maximised using gradient-based optimisation techniques. Since the KL on the RHS is between two multivariate Gaussians, there is an analytical expression for calculating this quantity, and estimate of the expectation can be obtained using numerical methods, in particular Monte Carlo sampling (together with the *reparametrisation trick* (Kingma and Welling 2014)).

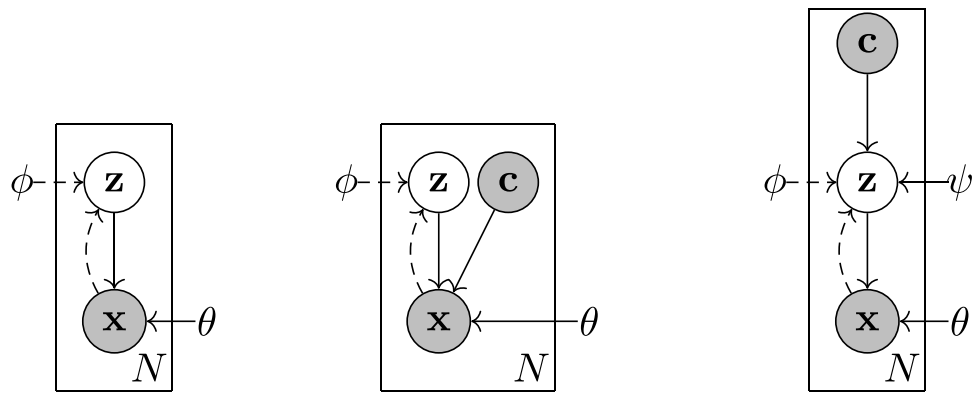
Are the latent representations induced by a VAE in any way *conceptual*? First, note that there is no pressure within the model to induce the sorts of factored representations in which the dimensions of \mathbf{Z} correspond to conceptual domains. Higgins et al. (2017) attempt to address this problem by introducing a weighting factor on the KL loss. Second, there is currently no mechanism in the model which allows concepts to be referred to using their names (e.g. *blue square*).

3.1.2 The conceptual VAE

One feature that we would like in the model is an explicit representation of the words or symbols that are used to refer to a concept (which we'll call the concept *label*). The obvious way to include the concept label in the model is as an explicit random variable \mathbf{c} . We could use a conditional VAE (Doersch 2016), with the label acting as an additional input into the decoder, so that when the decoder generates a data instance \mathbf{x} , it does so conditioned on \mathbf{c} as well as a point from the latent space \mathbf{z} (Fig. 2; centre). However, with this model there is no explicit representation of a concept (beyond its symbolic label). The key to the conceptual VAE is to introduce a new random variable for a concept label, \mathbf{c} , but introduce it at the very top of the graphical model (Fig. 2; right). The difference with the conditional VAE is that each concept \mathbf{c} now has an explicit set of parameters associated with it, which acts as \mathbf{c} 's representation.

⁹ In this section we use bold font for variables, e.g. the conceptual space \mathbf{Z} , to be consistent with the machine learning literature.

Fig. 2 Graphical models for the VAE (left), conditional VAE (centre) and the Conceptual VAE (right). Grey nodes represent observed variables and white nodes hidden variables



In terms of the generative story, first a concept label \mathbf{c} is generated, and then a point \mathbf{z} in the latent conceptual space is generated, *conditioned on* \mathbf{c} ; after that the generative story is the same as for the vanilla VAE: an instance \mathbf{x} is generated conditioned on \mathbf{z} . In this work we assume a uniform prior over the concept labels (more specifically a uniform prior over the atomic labels corresponding to each conceptual domain \mathbf{c}_i), and \mathbf{c} can effectively be thought of as a fixed input to the model, as provided by the data.

How do we model $p(\mathbf{z}|\mathbf{c})$? As before we use multivariate Gaussians with diagonal covariance matrices, but now the means and variances are *learnable parameters* ψ . We will sometimes refer to $p_\psi(\mathbf{z}|\mathbf{c})$ for a given concept \mathbf{c} as a *conceptual “prior”* (since these Gaussians replace the unit normal prior in the vanilla VAE), as well as \mathbf{c} 's learned representation. Since \mathbf{c} is factored, each \mathbf{c}_i has its own (univariate) Gaussian distribution; e.g., *red* will have its own mean and variance which define a Gaussian on the dimension corresponding to the COLOUR domain. It is this Gaussian which provides the answer to the question “what is the conceptual representation for *red*?”.

The ELBO equation now takes the following form:

$$\begin{aligned} & \log p(\mathbf{x}|\mathbf{c}) - \mathcal{D}(q(\mathbf{z}|\mathbf{x}), p(\mathbf{z}|\mathbf{x}, \mathbf{c})) \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}(q(\mathbf{z}|\mathbf{x}), p(\mathbf{z}|\mathbf{c})) \end{aligned} \tag{14}$$

How is this model trained, and what are the pressures that lead to conceptual representations being learned? For a training instance \mathbf{x} labelled with a concept \mathbf{c} , the training proceeds as before for the vanilla VAE: the encoder predicts a Gaussian $q(\mathbf{z}|\mathbf{x})$; this is sampled from (using the reparametrisation trick) to give a sample \mathbf{z}_s ; and $-\log p(\mathbf{x}|\mathbf{z}_s)$ is calculated to give the reconstruction loss. The key difference is in the calculation of the KL loss. Suppose that $\mathbf{c} = (\textit{green}, \textit{medium}, \textit{triangle}, \textit{bottom})$. The KL is calculated for each dimension, relative to the Gaussian for the particular atomic label for that dimension. For example, for the COLOUR domain (dimension

0), the KL would be between $q_\phi(\mathbf{z}_0|\mathbf{x})$ and $p_\psi(\mathbf{z}_0|\textit{green})$. So note that the supervision regarding the domains is provided here in the calculation of the KL.¹⁰ Unlike the vanilla VAE, the conceptual “priors” depend on the learned parameters ψ , which are the means and variances of the individual (univariate) Gaussians. We expect these learned means and variances to result in a neat separation along a dimension, since this will make it easier for the model to fit q to the conceptual representations, leading to a lower KL.

Conceptual space description Explicitly, in terms of our framework from Section 2.2, our conceptual model is given in the category **ConSp**, i.e. by a conceptual space. The model is $\mathbf{Z} = \mathbb{R}^n$, viewed as a product of n one-dimensional domains $\mathbf{Z} = \prod_{i=1}^n \mathbf{Z}_i$ with $\mathbf{Z}_i = \mathbb{R}$. An instance is a vector $\mathbf{z} = (z_1, \dots, z_n) \in \mathbf{Z}$. In particular for each image $\mathbf{x} \in X$ we obtain an instance via the (deterministic) encoder $q_\phi(\mathbf{x})$. Each concept label $\mathbf{c} = (\mathbf{c}_1, \dots, \mathbf{c}_n)$ defines a Gaussian fuzzy concept $\mathbf{c}(\mathbf{z}) = p(\mathbf{z} | \mathbf{c})$ with diagonal covariance matrix, as in Example 2. It forms a product concept over the domains as in Eq. 9, via:

$$\mathbf{c}(\mathbf{z}) = \prod_{i=1}^n \mathbf{c}_i(\mathbf{z}_i)$$

where each $\mathbf{c}_i(\mathbf{z}_i) = \mathbf{c}_i(\mathbf{z}_i; \mu_i, \sigma_i^2)$ is a one-dimensional Gaussian concept for concept label \mathbf{c}_i on \mathbf{Z}_i , with mean μ_i and variance σ_i^2 as trainable parameters.

3.1.3 A concept classifier

Here we show how the model can be adapted to act as a concept classifier. Note that, from a computer vision perspective, the classification task is trivial, and one that we would expect

¹⁰ The question of whether, and how, the level of supervision could be reduced and the domains learned automatically is an ongoing debate (Higgins et al. 2017; Locatello et al. 2019).

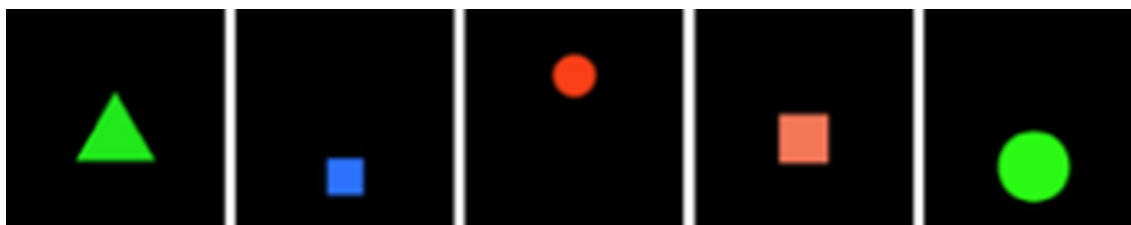


Fig. 3 Example shapes: (green, large, triangle, centre); (blue, small, square, bottom); (red, medium, circle, top); (red, medium, square, centre); (green, large, circle, bottom)

a well-trained CNN to solve. The classification task is being used here as a test of whether the induced conceptual representations can be employed in a useful way.

From a probabilistic perspective, the goal is to find the most probable concept c' given an input image x :

$$c' = \arg \max_c p(c|x) \tag{15}$$

$$= \arg \max_c p(x|c) \tag{16}$$

$$\approx \arg \max_c -\mathcal{D}(q(z|x), p(z|c)) + \text{recon_loss} \tag{17}$$

$$= \arg \max_c -\mathcal{D}(q(z|x), p(z|c)) \tag{18}$$

Line Eq. 16 follows from Eq. 15 because of the assumed uniform prior over concepts, and we use the ELBO from Eq. 14 as an approximation to the likelihood in going from Eqs. 16 to 17 (where recon_loss is the remaining part of the loss after the KL). The reconstruction loss is independent of c and so we end up with the satisfying form of the classifier in Eq. 18, in which the most likely concept for an input x is the one with the smallest KL relative to the encoding of x , as provided by q .

3.2 Experiments

We use the Spriteworld software (Watters et al. 2019) to generate simple images. These consist of coloured shapes of particular sizes in particular positions in a 2D box, against a black background. For the main dataset, there are three shapes: {square, triangle, circle}; three colours: {red, green, blue}; three sizes: {small, medium, large}; and three (vertical) positions: {bottom, centre, top} (see Fig. 3). We ran the sampler to generate a training set of 3,000 instances, and development and test sets with 300 instances each. Appendix A contains the parameters used in the Spriteworld software to generate the main dataset.

The encoder, which takes an image x as input, is instantiated as a CNN, with 4 convolutional layers followed by a fully-connected layer. A final layer predicts the means and variances of the multivariate Gaussian $q_\phi(z|x)$. The ReLU activation function is used throughout. The decoder, which

takes a latent point z as input, is instantiated as a de-CNN, with essentially the mirrored architecture of the encoder. The reconstruction loss we use on the decoder for predicting the pixel values in an image x is the MSE loss.

The implementation was in Tensorflow. The full set of parameters to be learned is $\theta \cup \phi \cup \psi$, where θ is the set of parameters in the encoder, ϕ the parameters in the decoder, and ψ the means and variances for the conceptual representations (12 each for the main dataset). The training was run for 200 epochs (unless stated otherwise), with a batch size of 32, and the Adam optimizer was used. Finally, we added 2 “slack” dimensions to the latent space Z , in addition to the 4 dimensions for the conceptual domains. These slack dimensions are intended to capture any remaining variability in the images, beyond that contained in the concepts themselves. Appendix B contains more details of the neural architectures used in our experiments, including the various hyper-parameter choices.

3.2.1 Clustering effects and classification accuracy

Figure 4 shows the means and log-variances predicted by the encoder for each dimension, for a set of instances, with the colour-coding indicating the atomic concept labels from the different domains. For example, in the set of 4 plots at the top left, the means and log-variances for dimension 0 are plotted; and in the top-left of those 4 plots, each point is colour-coded with the colour of the corresponding instance. What this plot shows is the neat separation for the means along the COLOUR dimension, for each of the 3 colours. The other 3 plots contain the same set of points, but colour-coded with atomic labels from the remaining domains of SIZE, SHAPE and POSITION. With the 3 remaining plots we expect to see no discerning pattern, since we would like the first dimension to encode COLOUR only (although note that, in this particular training run, dimension 1—corresponding to SIZE—does appear to be encoding some information about the colour).

The plots were created using the model evaluated on classification accuracy below, which performed well on the development data. The instances were taken from the training

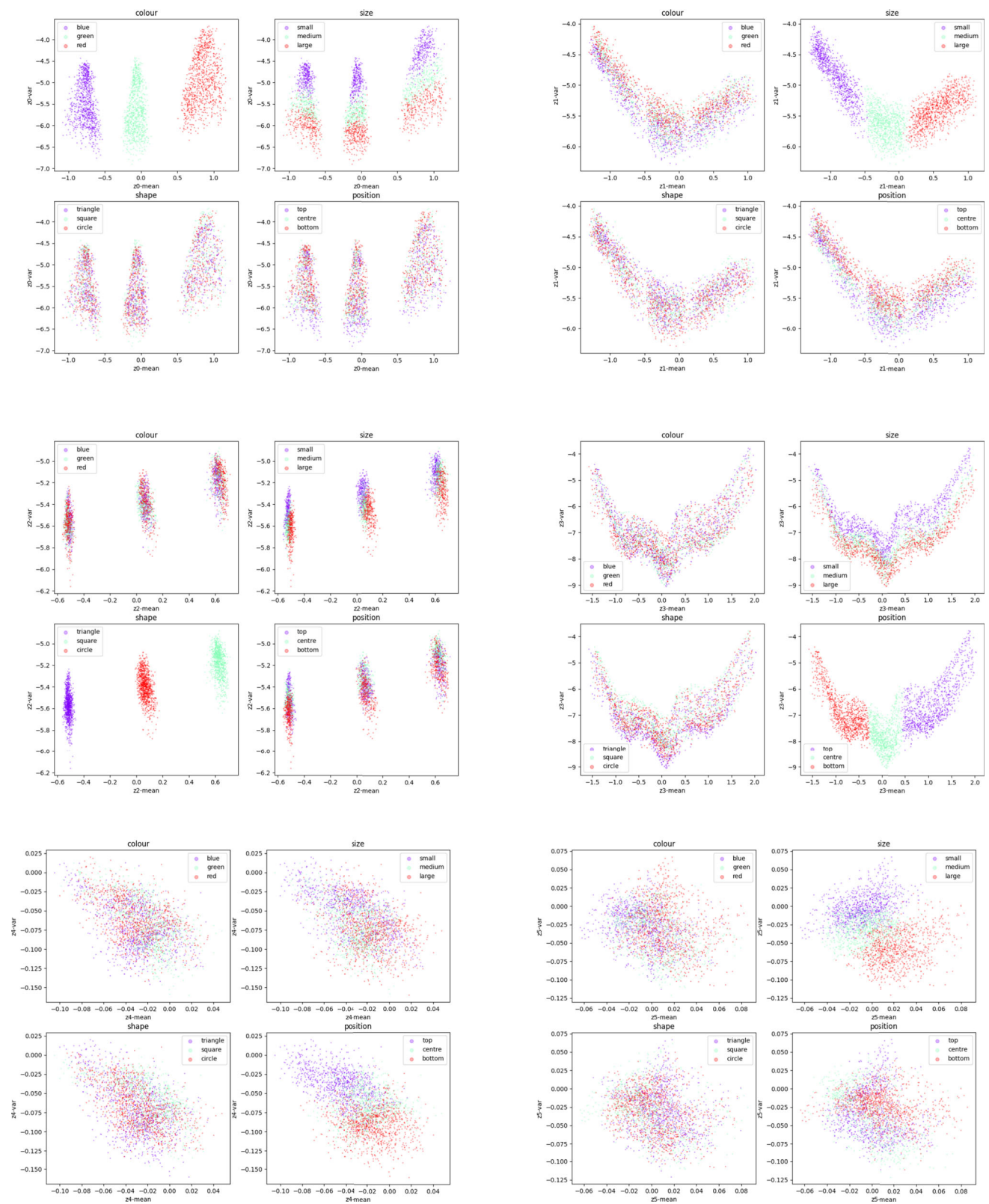


Fig. 4 The means and log-variances for each dimension predicted by the encoder, for a set of instances; means on the x-axis, log-variances on the y-axis. Colour-coding, from top-left clockwise: COLOUR, SIZE, POSITION, SLACK- DIM- 2, SLACK- DIM- 1, SHAPE

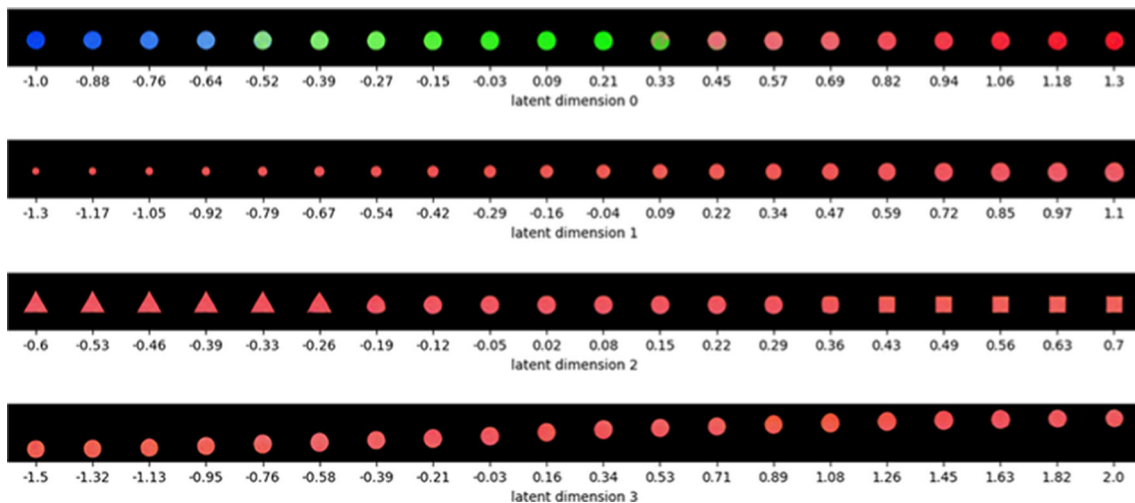


Fig. 5 Traversals along each latent dimension for a large red circle in the centre

data.¹¹ The plots in the top-right are for dimension 1 (corresponding to SIZE), and again we obtain a neat separation for the means, when colour-coded with the size of the instance, with instances labelled *medium* sitting in the middle. The second-row (right) plots are for dimension 3 (POSITION), and again we see a neat separation of the means with instances labelled *centre* sitting between those labelled *top* and *bottom*. The second-row (left) plots are for dimension 2 (SHAPE). Here we see a clear separation with the predicted means occupying a short range, which reflects the discrete nature of these concepts. The plots for the slack dimensions are in the bottom row, with no discernible pattern (except perhaps in the SIZE dimension bottom-right).

We evaluated the same model as a classifier, using the formulation in Eqs. 15, 16, 17 and 18 above. The accuracy on the development data for the COLOUR and SHAPE domains was 100%, with accuracies above 98% for the other two domains. These high accuracies transferred over to the test data.

3.2.2 Continuity within domains

Figures 5 and 6 provide further qualitative demonstration of how the conceptual domains are neatly represented on each dimension. An instance of a large red circle in the centre and a medium-sized blue square at the bottom are passed through the encoder, giving a mean for each of the 4 dimensions. Then, the mean value is systematically varied for one of the dimensions only (through regular increases and decreases), keeping the other 3 fixed. All resulting combinations of the 4 mean values are then input to the decoder, giving the images

in the figure.¹² What the transitions clearly demonstrate is not only how one latent dimension encodes just one domain, but also how the concepts smoothly vary along one dimension. Note how in both examples dimension 2 encodes a shape somewhere between a *triangle* and a *circle*, and also a shape somewhere between a *circle* and a *square*. Dimension 1 shows a smooth transition from *small* to *medium* to *large*, and dimension 3 from *bottom* to *center* to *top*.

In order to investigate these ordering effects further, we created a new dataset which contains all the colours of the rainbow, with the same shapes, sizes and positions. The continuous colour ranges now cover a much larger proportion of the range of possible values (see Appendix A.1), with the occasional gap (e.g. between *green* and *blue*). The training data again consisted of 3,000 randomly generated instances, with a development set of 300 instances.

Again we chose a trained model which performed well on the classification task on the development data (with accuracies well into the 90s for all domains), and plotted the colour-coded means and variances as predicted by the encoder. Figure 7 again shows a neat separation for the COLOUR domain, with very similar patterns for the other domains (not shown), and to those exhibited in Fig. 4. Looking carefully at the plot in the top-left, we see that the colours are not only neatly separated along the COLOUR dimension, but also that the ordering of the rainbow is faithfully represented: *blue*, *indigo*, *violet*, *red*, *orange*, *yellow*, *green*. Figure 8 shows an example traversal along the COLOUR dimension only, for the colour-extended dataset, again demonstrating an ordering consistent with a rainbow.

¹¹ The same patterns were observed on the development data. We used the training data since this gives denser plots.

¹² The idea of plotting transitions along a dimension is taken from Higgins et al. (2017).

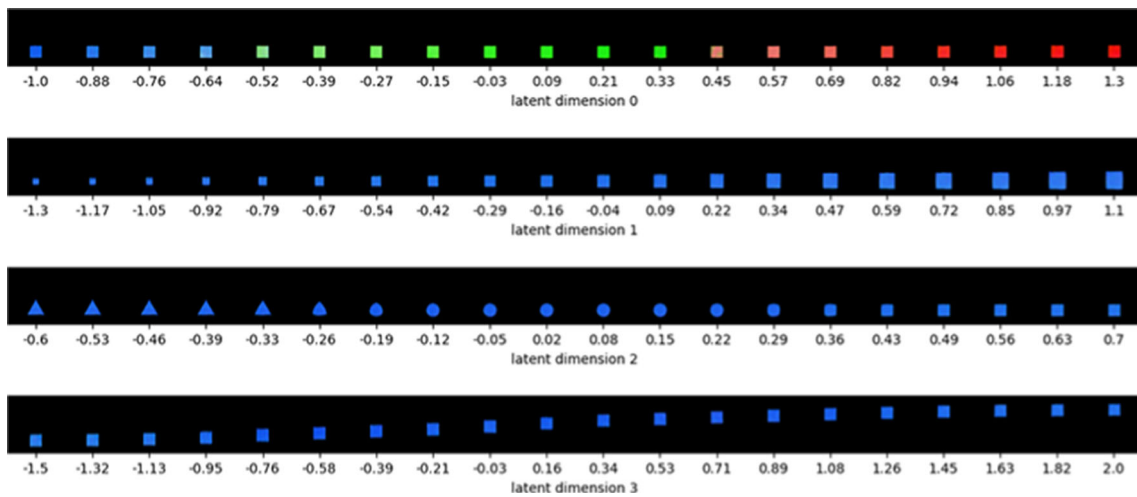


Fig. 6 Traversals along each latent dimension for a medium-sized blue square at the bottom

Fig. 7 The means and log-variances for COLOUR predicted by the encoder, for the rainbow colour set

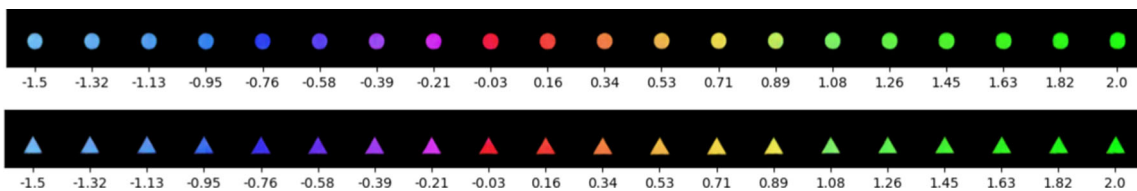
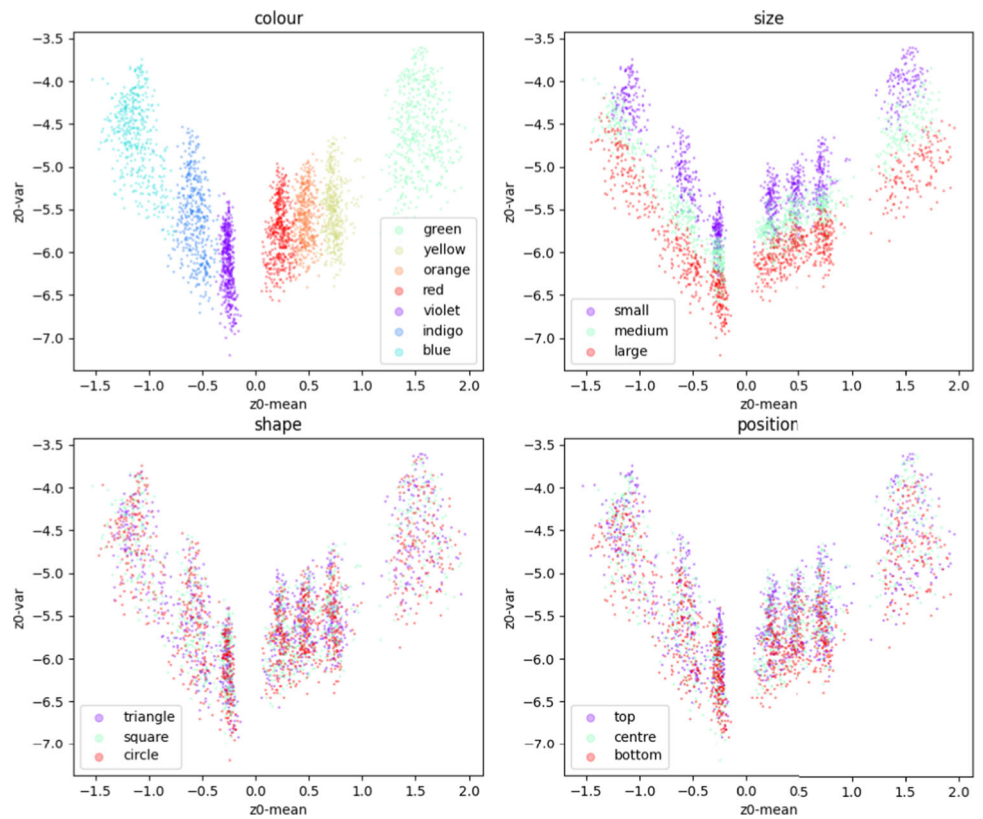


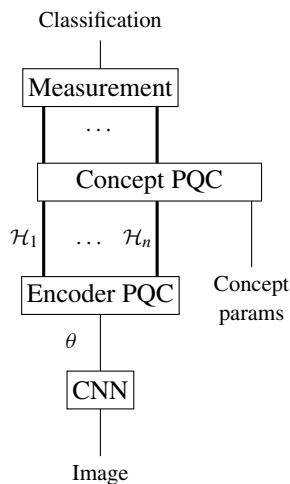
Fig. 8 Traversals along the COLOUR dimension for two examples from the colour-extended dataset

4 Quantum implementation: a hybrid network with PQC

In this section we also set up a probabilistic learning objective in order to induce conceptual representations, but using a discriminative classifier rather than a generative model. In addition, the classifier is implemented as a hybrid network consisting of a classical convolutional neural network (CNN) (Goodfellow et al. 2016, Ch.9) followed by a parameterised quantum circuit (PQC) (Benedetti et al. 2019). We use the network to classify the same set of images and labels from the classical experiments in Section 3.

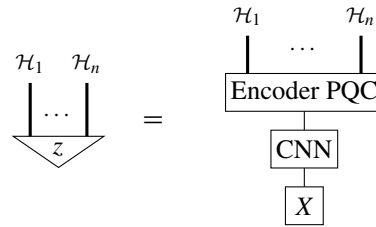
4.1 The hybrid network

An input image is first processed by a CNN which outputs classical parameters which are fed into a PQC. This PQC we call the *encoder PQC*; it implements a quantum state z which is the representation of the image in our model. Given a concept C , a separate *concept PQC* implements a quantum effect corresponding to C which can be applied to the instance z , as described in Sections 2.2 and 2.4. We assume that the factorisation of the model into the domains $\mathcal{H}_1, \dots, \mathcal{H}_n$ is known by the model; in our experiments these will be the four domains SHAPE, COLOUR, SIZE, POSITION. The overall setup is shown below, with thin wires denoting classical data and each thick wire denoting a Hilbert space given by some number of qubits.

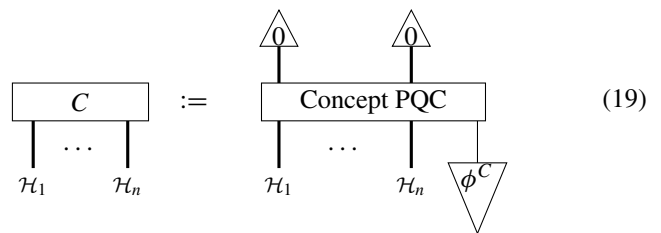


Given an input image and the parameters encoding a concept, a single run of the circuit produces a “yes” or “no” to determine whether the concept has been deemed to fit the image. The probability of each outcome is obtained either by sampling the circuit many times (on a physical device) or direct calculation (in simulation). With the probabilities for each concept, one can then classify which concept best fits the input image.

In more detail, each instance z is a pure quantum state given by passing an image X into the CNN and then using the resulting parameters in the encoder PQC network:



Each specific concept C can be understood as a measurement with two outcomes “yes” and “no”, such that outcome “yes” means the concept has been deemed to fit the instance. The measurement is given by a Pauli-Z measurement on each qubit, with the overall outcome “yes” identified with obtaining outcome 0 on every qubit individually, and all other outcomes labelled as “no”. Diagrammatically this is expressed as follows:

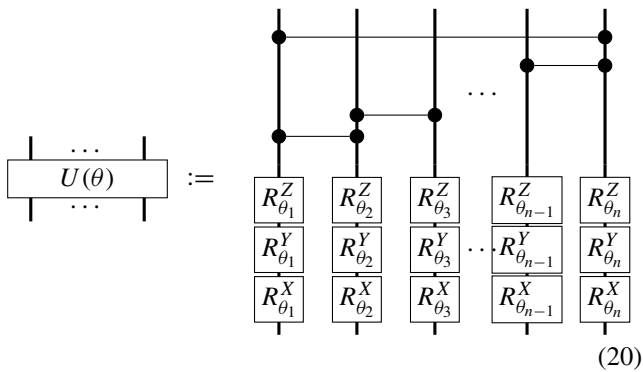


where ϕ^C are the parameters encoding the concept C . Each concept C can be either pure or mixed, depending on whether a pure or mixed circuit is chosen for the concept PQC, which we discuss in Section 4.1.1.

4.1.1 The CNN and PQCs

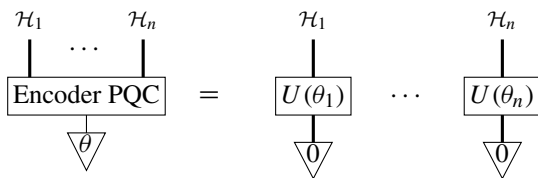
We use the same CNN from the classical experiments in Section 3.2 for the image processing. For the classical experiments the CNN predicted the means and variances of a multivariate Gaussian, whereas here the CNN predicts the parameters of the encoder PQC. The PQCs make use of the parameterised circuit ansatz shown below, defined over any finite collection of qubits. The ansatz $U(\theta)$ is given by performing parameterised X, Y, Z rotations on each qubit, followed by entangling pairs of adjacent qubits using controlled Z gates (with an additional gate operating on the two outermost qubits to complete the chain). Multiple layers of this ansatz can be composed to give a more complex circuit. We define another ansatz $V(\theta)$ in the same way but with initial rotations in the reverse order Z, Y, X . An important special case is that, when given on a single qubit, $U(\theta)$ is simply equal to sequential parameterised X, Y and Z

rotations. Similarly $V(\theta)$ on a single qubit amounts to rotating in the order Z, Y, X .



In the above each $\theta_j = (\theta_{j,X}, \theta_{j,Y}, \theta_{j,Z})$ consisting of three single parameters passed respectively to the X, Y, Z rotations on qubit $j = 1, \dots, n$. All are in turn contained in the parameters vector θ . In fact this ansatz is universal in that with sufficient layers of the form $U(\theta)$ one may implement any unitary circuit.¹³

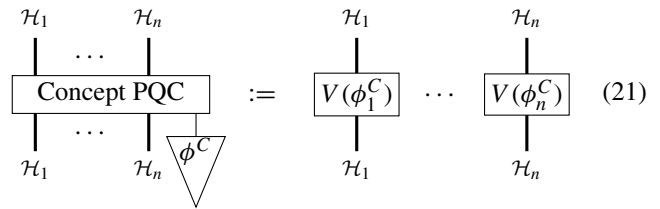
Now let us describe the encoder and concept PQCs in more detail. Both consist of some number of qubits per domain \mathcal{H}_i . The form of the encoder PQC is as follows:



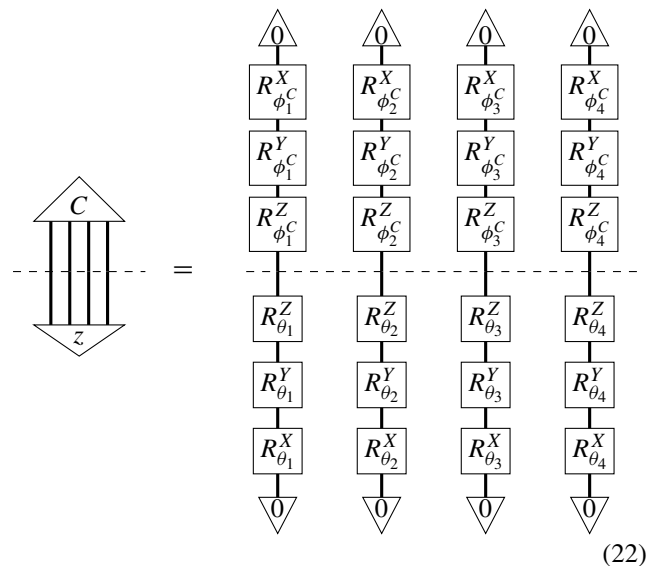
More generally we can compose multiple layers of such U circuits on each domain. Here the $|0\rangle$ states denote product states $|0 \dots 0\rangle$ on each \mathcal{H}_i . Thus by construction the encoder never involves entanglement across domains, and can be viewed as a single encoder per domain. Since the ansatz U is universal, the encoder is able to prepare an arbitrary quantum instance.

In the initial basic setup used, beginning in Section 4.2, we only have one qubit per domain \mathcal{H}_i , and only use one layer in the encoder. In this case the encoder simply carries an X, Y and Z rotation per qubit, involving no entanglement. In this

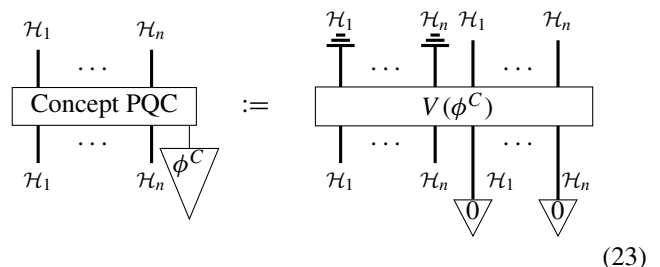
basic setup, the concept PQC also involves no entanglement, taking the following form.



Concretely, with four domains and one qubit per domain, in this setup the application of a concept C to an instance z amounts to the (probability of the) circuit shown below with post-selection, where θ is the encoding of the image from the CNN, ϕ^C are the learned concept parameters and each wire is a single qubit.



In order to capture mixed and entangled concepts, in Section 4.4 we use a richer form for the concept PQC. Entanglement is provided by using the full ansatz $V(\theta)$ over all domains. To introduce mixing, we use an ancillary copy of each domain $\mathcal{H}_1, \dots, \mathcal{H}_n$, prepared in initial state $|0\rangle$, and then discard the original domains as in the following circuit:



More generally one can include multiple V layers prior to discarding. Note that since this ansatz is universal we can

¹³ The entangling layer is self-inverse, so that two layers allow us to implement a rotation on any qubit. A swap operation on any pair of qubits can be implemented using three layers, and from this any CX gate. Hence we may implement the universal gate set given by single-qubit phase and Clifford gates; see, for example, Van de Wetering (2021).

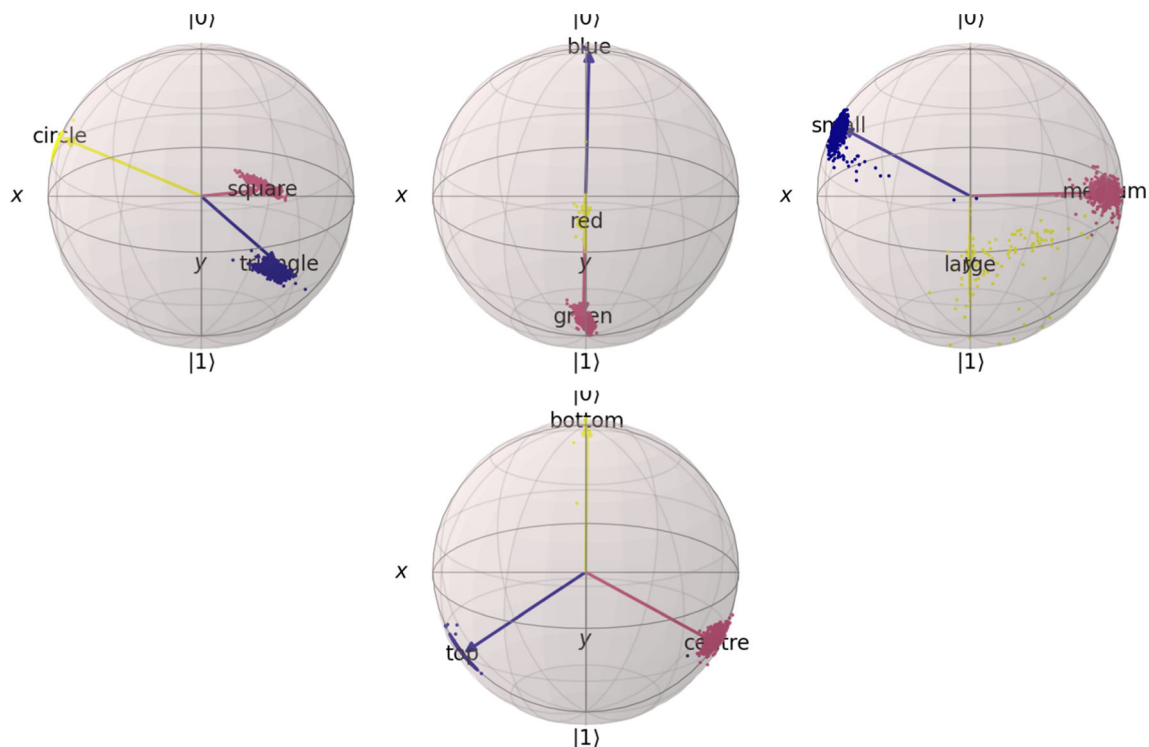


Fig. 9 Visualisation of the pure concept effects and instance states on the Bloch sphere, for SHAPE, COLOUR, SIZE and POSITION

implement any unitary with sufficient layers of the form V , and thus any (sub-normal) quantum concept.

4.1.2 Discriminative training

The classical concepts model from earlier is a generative model consisting of an encoder and a decoder. Here we choose to train the quantum model to perform binary classification; hence the basic model is a discriminative model with an encoder only.¹⁴ The loss function is the standard binary cross entropy (BCE) loss for binary classification. The full set of parameters to be learned is $\psi \cup \phi$, where ψ is the set of parameters in the classical encoder CNN and ϕ is the set of PQC parameters associated with the set of 12 basic concepts.

The training data contains the 3,000 positive examples from Section 3.2 and an additional 3,000 negative examples. Each negative example is created from a positive one by randomly sampling an incorrect concept for each domain; for example, if the positive example is *(green, large, triangle, centre)* then a negative example could be *(blue, medium, square, bottom)*. Since we are effectively learning each domain independently in the basic model, a negative example disagrees on every domain. Later models will use variations on this data.

¹⁴ In Section 4.2.2 below we investigate how the addition of a decoder can affect the instance and concept representations.

The implementation is in Tensorflow Quantum, and the whole hybrid network—both the quantum and the classical parts—are trained end-to-end in simulation on a GPU. The training was run for 100 epochs (unless stated otherwise), with a batch size of 64 (32 images, each with a positive and negative label), and the Adam optimizer was used.

4.2 Instance states and concept effects

We trained a quantum model, using the circuit shown in Eq. 22 above, and tested it on the 300 examples in the development set. The model was trained to perform binary classification, but at test time we choose the concept for each domain which has the highest probability of applying to the input image. The classification model performed with almost perfect accuracy, obtaining 100% on the COLOUR and SHAPE domains, and 99% and 97% on the POSITION and SIZE domains, respectively. This high accuracy carried over to the 300 examples in the test set, obtaining 100% on the COLOUR and SHAPE domains, and 96% and 97% on the POSITION and SIZE domains, respectively.

Figure 9 visualises the pure effects for each of the 3 concepts on the 4 domains, by plotting the corresponding pure states on a Bloch sphere. We are able to perform the visualisation for this basic model since only one qubit is being used per domain, with no entanglement. The clusters of dots around each concept are the corresponding instances (pure

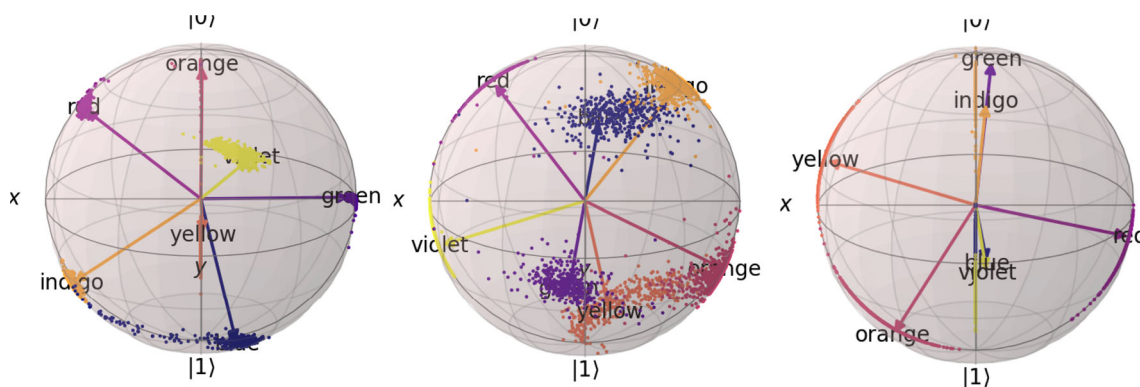


Fig. 10 Visualisation of the concept effects and instance states on the Bloch sphere, for 3 trained models, for the COLOUR domain on the rainbow dataset

states) in the training data. This visualisation is for the model which performs as described above on the classification task; a model trained from a different random initialisation would have the concepts and instances distributed differently around the sphere, but this visualisation is representative in terms of how the concepts are typically separated and the instances clustered. Note how the 3 concepts on each domain are being pushed apart (strikingly so in the case of the POSITION domain) and how the concepts sit neatly in the centre of each cluster of instances.

4.2.1 The rainbow dataset

In order to test our model further, we used the rainbow dataset from Section 3.2, and in order to train the discriminative model, we added a further 3,000 negative examples (for each epoch) to the 3,000 positive ones, randomly generated as before. Perhaps unsurprisingly, it was more difficult with this data to obtain a clean separation of the colours on a single qubit.¹⁵ However, with a weighting of 0.5 applied to the negative examples in the binary cross-entropy loss, and running the training for 200 epochs, we were able to obtain the distribution of colours around the Bloch sphere shown in Fig. 10 (with instances again taken from the training data). The three visualisations are for three separately trained models (i.e. with three different random initialisations of the model parameters).

In terms of accuracy on the development data, the classification model for the Bloch sphere on the far left achieved similar scores on the non-colour domains as before, and an overall accuracy of 95% on COLOUR, with F1-scores ranging from 91% to 100% for the individual colours. The Bloch sphere in the middle is for a model with similar performance,

and is shown to demonstrate the variation in models. The example on the far right is cherry-picked as an example of how the training is able to neatly represent the various colours on the Bloch sphere: note how the *yellow*, *orange* and *red* instances are beautifully placed on the circumference of a circle, with the *red* instances leading into *orange* and then *yellow*.

4.2.2 Adding a decoder loss

One notable feature of the visualisations in Fig. 9 is how “tight” the instance clusters are, despite the variation in the images for a single concept (for example the variation in red shapes in Fig. 3). There may be use-cases where we would like the representation of instances to better reflect the variation in the underlying images, for example in order to better capture correlations across domains (see Section 4.3 below).

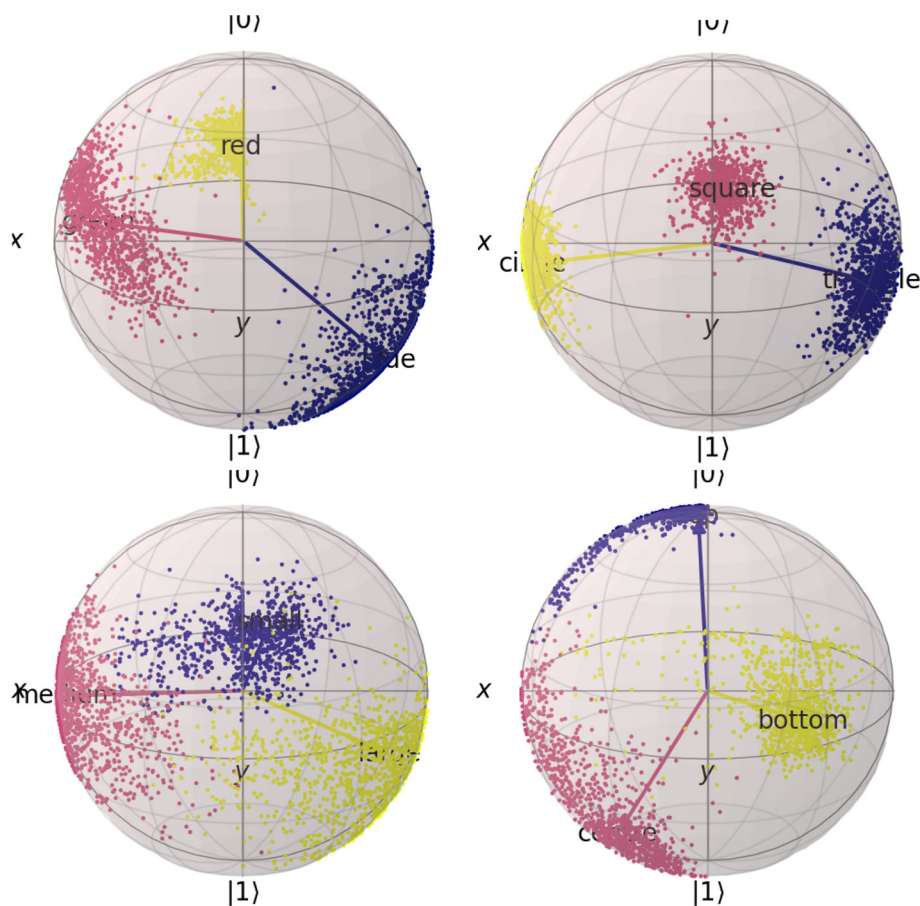
In order to provide more of a “spread” of the instances, we experimented with an additional decoder loss in the loss function below, where BCE is the binary cross-entropy loss, D is the data with N instances $\{X_i\}_i$, and ψ and ϕ are the parameters of the encoder network:

$$\text{Loss}(D, \psi, \phi, \chi) = \text{BCE}(D, \psi, \phi) + \frac{\lambda}{N} \sum_i \text{SE}(\text{DeCNN}(\chi, \text{CNN}(\psi, X_i)), X_i) \tag{24}$$

The decoder is a deconvolutional neural network (DeCNN), with parameters χ , which essentially is the CNN “in reverse”: it takes as input the angles output by the CNN, given an image X_i , and outputs RGB values for each pixel in the image. SE is the sum of squared errors across all RGB values in the image, and λ is a weighting term in the overall loss. The intuition is that, in order to obtain a low SE loss, the encoder CNN has to output angles which are sufficiently informative in

¹⁵ Of course there is nothing to prevent us from using more than one qubit per domain, in order to provide a larger Hilbert space in which to represent the additional colours, but the visualisation is harder with more qubits.

Fig. 11 Visualisation of the concept effects and instance states for all 4 domains, for the basic dataset with an additional decoder loss



order for the DeCNN to accurately reconstruct the original image. Now the model is similar to the Conceptual VAE (albeit without the generative model interpretation), in that it has both “encoder” and “decoder” parts to the loss.¹⁶

Figure 11 shows how the instances can be distributed more broadly around the Bloch sphere, using the additional decoder loss (with $\lambda = 0.1$). This model still performs well as a classification model on the development data, achieving 98% accuracy on SIZE, 99% on COLOUR, 100% on SHAPE, and 98% on POSITION. As a qualitative demonstration of this approach, note how the instances for *centre* and *top* start to merge into each other (blue and red instance dots bottom right), and also for *medium* and *small* (blue and red instance dots bottom left), which is what we would expect for a less discrete representation.

4.3 Capturing correlations

Here we show how one of the characteristic features of quantum theory, namely entanglement, can be used to capture correlations across domains. In order to test whether our

model can handle concepts which contain correlations, we define a new concept which we call *twike*, which is defined as (*red and circle*) or (*blue and square*) (i.e. it applies to images containing red circles or blue squares). Figure 12 shows some examples of *twikes* and non-*twikes*.

The concept PQCs we have considered so far, of the form in Eq. 21, are unable to learn the concept *twike*, since the domains have been treated independently, with each of the 4 domains effectively containing its own independent concept. In order to create connections between the domains in the concept PQC, we can apply our full ansatz V from Section 4.1.1, involving controlled-Z gates between wires, across multiple domains. In this first experiment we assume knowledge of the fact that, for the *twike* concept, the correlations are across the SHAPE and COLOUR domains, with entangling gates only between the qubits for SHAPE and COLOUR. (This assumption will be relaxed for some of the experiments below.) We also assume that the remaining domains are not relevant and so are not measured, thus effectively being discarded in the concept. We apply potentially multiple layers of ansatz V to the relevant domains, and so the resulting form of the *twike* concept over the four domains is as shown in Fig. 13, where ϕ are the learned parameters for the *twike* concept.

¹⁶ One possibility for future work is to develop and implement a “quantum VAE” (Khoshaman et al. 2018) for concept modelling, and have a generative model in which all parts of the model are quantum.

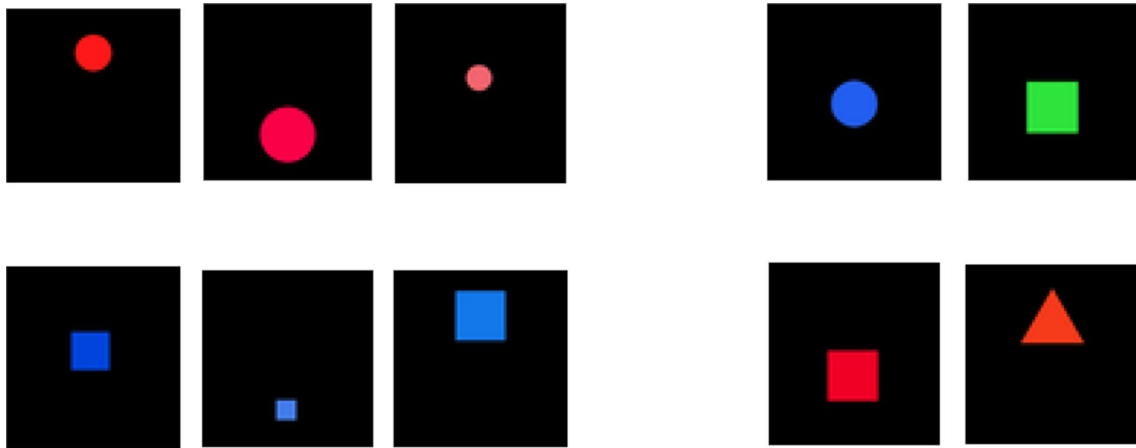


Fig. 12 Examples of twikes (on the left) and non-twikes (on the right)

The training of this model only updates the rotation parameters of the concept PQC; the parameters of the encoder (i.e. the CNN) are kept fixed from the earlier training of the basic model. The loss function is binary cross entropy, as before, with the 3,000 examples from Section 3.2 used as training data. Roughly 20% of these instances are positive examples of *twike*, with the remaining being negative examples. We trained this model for 50 epochs, using 2 layers of the rotation and entangling V ansatz for the concept PQC, and obtained 100% accuracy on the unseen test examples. It was only through the introduction of the entangling gates that we were able to learn the *twike* concept at all.

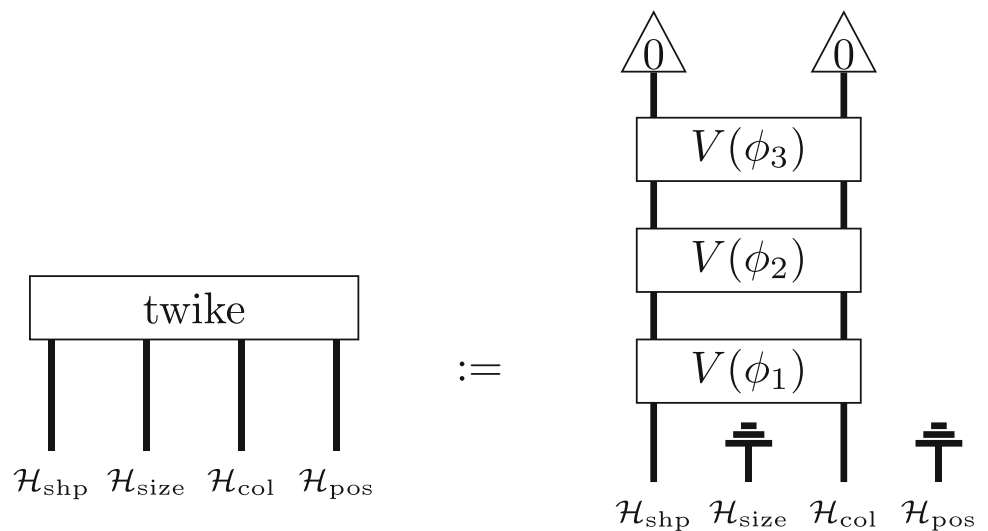
In terms of the discussion of entanglement and classical correlation in Section 2.5, we can say that the *twike* concept can be naturally described without entanglement, as a classical combination of the pure concepts *red circle* and *blue square* (at least in the case where these pure effects are orthogonal). However, such correlations are not always

immediately implementable in many conventional classical models. In terms of the Conceptual VAE, it would be possible to capture correlations using the covariance matrix of the multivariate Gaussian. However, a standard assumption in VAEs is to assume a multivariate Gaussian with a diagonal covariance matrix (and so no correlations across domains). Whether a concept like *twike* could be easily modelled using the Conceptual VAE, especially as the number of domains is increased, is left as a question for future work.

4.4 Learning general mixed and entangled concepts

One assumption made above in the *twike* experiments was that the relevant domains—in this case SHAPE and COLOUR—are known in advance, so that the concept PQC can effectively ignore the wires corresponding to the other domains. One interesting question is whether the concept PQC could also

Fig. 13 Encoder PQC for learning *twike*, here shown with 3 layers of the rotation and entangling V ansatz



learn which domains are relevant, as well as which of those domains should be correlated, if provided with all 4 wires as input. To allow for such correlations between arbitrary domains, the concept PQC should allow for entanglement between any of its domains. Furthermore, to enable discarding of domains, we require mixed quantum effects. Both of these features can be included by using our most general form of the concept PQC Eq. 23.

In order to test the learning of these general concepts, we set up a similar experiment to *twike*, but this time with just *red* as the concept to be learned. Of course the encoder had already learned *red* when trained to perform classification in the basic setup, but in this experiment we remove knowledge of which wire the COLOUR domain is on, and see whether a new concept PQC can learn *red*, given red and non-red instances as input.

Again the training of this model only updates the rotation parameters of the concept PQC; the parameters of the CNN are kept fixed. The loss function is again binary cross entropy, with the usual 3,000 examples as training data. Roughly 33% of these instances are positive examples of *red*, with the remaining being negative examples. We trained this model for 50 epochs, using 2 layers of rotation and entangling gates for the concept PQC, and obtained 100% accuracy on the unseen test examples. It was only through the introduction of the discarding (plus entangling gates) that we were able to obtain these high accuracies.

4.5 Concepts containing logical operators

For one final set of experiments, we investigated whether the entangling and discarding PQC Eq. 23 could learn concepts built from logical operators, with concepts such as *red or blue*.

4.5.1 Conjunction across domains

The first concept with a logical operator that we consider is *red and circle*, firstly with the knowledge of which domains are relevant for the concept (in this case COLOUR and SHAPE). The encoder PQC is the simple one from Eq. 22, but with only the COLOUR and SHAPE wires (so the other two are effectively discarded). We used the same 3,000 training examples, of which roughly 17% are positive examples and 83% negative examples. In this case the learning is particularly easy, and the model obtains 100% accuracy with only a single layer of rotations for the PQC, without any entangling gates or discarding of any ancilliary qubits. The reason is that the factorisation of the domains through the tensor product has effectively provided all the structure required to use conjunction.

When the knowledge of which domains are relevant is removed, and the more general encoder PQC in Eq. 23 is

used, learning becomes harder but an encoder PQC with 4 layers of rotation and entangling gates is able to learn the concept with 100% accuracy.

4.5.2 Disjunction within domains

Next we consider disjunction, but *within* rather than *across* domains, with the concept to be learned being *red or blue*. Of the 3,000 training examples, 61% are positive examples and 39% negative. Again, when knowledge of which domains are relevant is provided to the concept PQC, the learning is easy, with 100% accuracy obtained with a single layer of rotations.

If each point on the Bloch sphere were to correspond to an instance of the COLOUR domain, i.e. a single colour, as in our model, then the PQC learning such a pure effect for *red or blue* will in fact be simply learning a single colour, intuitively somewhere “in between” *red* and *blue*. When the domain only comes with a few concepts, such as the 3 concepts used here, this single instance may do well in approximating *red or blue*, as with the 100% accuracy. However, in the presence of more concepts, we expect that a concept for *red or blue* should involve mixing. And when knowledge of which domains are relevant is not provided to the PQC, *red or blue* can also be successfully learned with the more general PQC in Eq. 23 with 3 layers of rotation and entangling gates, including discarding.

5 Related work

The Conceptual VAE is inspired by Higgins et al. (2017), who introduced the β -VAE for unsupervised concept learning. However, the focus of Higgins et al. (2017) is on learning the conceptual *domains*, i.e. the underlying factors generating the data (Bengio et al. 2013), which they refer to as learning a *disentangled* representation. The main innovation to encourage the VAE to learn a disentangled, or factored, latent space is the introduction of a weighting term β on the KL loss. Higgins et al. (2017) show that setting β to a value greater than 1 can result in the dimensions of \mathbf{Z} corresponding to domains such as the lighting or elevation of a face in the celebA images dataset, or the width of a chair in a dataset of chair images. Our focus is more on the conceptual representations themselves, assuming the domains are already known, and the question of how concept labels can be introduced into the VAE model.

A paper in NLP that uses a model very similar to the Conceptual VAE is Bražinskas et al. (2018) which introduces the Bayesian skip-gram model for learning word embeddings. One key difference which distinguishes our work from the word embeddings typically used in NLP is that we do not restrict ourselves to the textual domain, meaning that our conceptual representations are *grounded* in some other modality

(in our case images) (Harnad 1990), bringing them closer to the human conceptual system. Another relevant paper from the NLP literature, which does consider grounding, is Schlangen et al. (2016), where the meanings of words are treated as classifiers of perceptual contexts, similar to how we use classification to induce conceptual representations.

The Conceptual VAE uses Gaussians to represent concepts, since they are the typical distributions used with VAEs and because they are convenient from a mathematical perspective. However, the use of Gaussians is also prevalent in the neuroscience literature, appearing for example as the *Laplace assumption* in the “free-energy” or “predictive processing” framework (Friston and Kiebel 2009; Bogacz 2017).

In terms of the quantum models, Smolensky has a large body of work arguing for tensor product representations in linguistics and cognitive science more broadly (Smolensky and Legendre 2006). Recently these techniques have been integrated into neural models for NLP (Huang et al. 2018). Another line of work which associates tensor-product representations with grammatical structure is the “DisCoCat” research program attempting to build distributed, compositional representations of language, which began with Coecke et al. (2010). Recently this work has culminated in the running of quantum NLP models on real quantum hardware (Lorenz et al. 2023).

The field of *quantum cognition* (Pothos and Busemeyer 2013) has already been mentioned. Some recent work in this area includes Epping and Busemeyer (2022) and Epping et al. (2021), where the latter is concerned with modelling human judgements of colour similarity and uses a Hilbert space representation similar to our models. The learning of concepts containing logical operators has a formal connection to quantum logic (Birkhoff and von Neumann 1936) and Boolean concept learning in general, for which there is a large literature (Goodwin and Johnson-Laird 2013).

6 Conclusion and further work

In this article we have presented a new modelling framework for structured concepts using a category-theoretic generalisation of Gärdenfors’ conceptual spaces, and shown how the conceptual representations can be learned automatically from data, using two very different instantiations: one classical and one quantum. The main contributions of this foundational work are the category-theoretic formalisation, and the two practical demonstrations, especially the quantum implementation which is particularly novel. Substantial further work is required to demonstrate that the framework can be applied fruitfully to data from a psychology lab, which would connect our work directly with quantum cognition, and also to agents acting in (virtual) environments, which would connect it to agent-based AI (Abramson et al. 2020).

In future more interpretative work on quantum concepts is needed to clarify their advantages, such as those offered by entanglement discussed in Section 2.5, and their naturality as a model in cognition. Another benefit of quantum models over conceptual spaces not explored here is the presence of a “negation” C^\perp on concepts with $C \leq \ddagger$ (Rodatz et al. 2021; Shaikh et al. 2021). In contrast, negation is harder to define for concepts in conceptual spaces; for example the complement of a convex region is generally non-convex.

Another interesting question is whether the Conceptual VAE can be applied to data generated from a conceptual hierarchy—for example having shades of colour such as *dark-red*—and whether the learned Gaussian representations for concepts can be partially ordered in an appropriate way (Clark et al. 2021). The quantum concepts as effects have a natural ordering, as discussed in Section 2.4, and it would be an interesting comparison to see if hierarchies could be more easily learned with the quantum models.

To make full use of the compositional approach, one should also describe conceptual *processes*, such as reasoning processes and “metaphorical” mappings between domains, now given by CP maps between quantum models. One could then compare these with the processes in the category **ConSp** of fuzzy conceptual processes from Tull (2021).

Finally, even though all the practical work here has been carried out in simulation on a classical computer, the number of qubits is relatively small, and the circuits relatively shallow, and so the running of these models on real quantum hardware is a distinct possibility. Also left for future work is the search for tasks which could demonstrate advantages for our quantum representations, for example establishing whether non-separable effects in the theory do provide an advantage over classical correlation in modelling conceptual structure.

A The shapes dataset

The parameters used in the Spriteworld software to generate the Shapes dataset:

```
COLOURS = {
  'red': distribs.Mixture([distribs.Continuous('c0', 0.95, 1.), distribs.Continuous('c0', 0., 0.95)]),
  'blue': distribs.Continuous('c0', 0.55, 0.65),
  'green': distribs.Continuous('c0', 0.27, 0.37),
}

SHAPES = {
  'triangle': shapes.polygon(num_sides=3, theta_0=np.pi/2),
  'square': shapes.polygon(num_sides=4, theta_0=np.pi/4),
  'circle': shapes.polygon(num_sides=30),
}

SIZE = {
  'small': distribs.Continuous('scale', 0.1, 0.17),
  'medium': distribs.Continuous('scale', 0.17, 0.23),
  'large': distribs.Continuous('scale', 0.23, 0.3),
}

POSITION = {
  'top': distribs.Product([distribs.Continuous('y', 0.56, 0.74), distribs.Discrete('x', [0.5])]),
  'centre': distribs.Product([distribs.Continuous('y', 0.38, 0.56), distribs.Discrete('x', [0.5])]),
  'bottom': distribs.Product([distribs.Continuous('y', 0.2, 0.38), distribs.Discrete('x', [0.5])]),
}
```


Additional parameters for the COLOUR domain:

```
distributions.Continuous('c1', 0.5, 1.), #saturation
distributions.Continuous('c2', 0.9, 1.), #brightness
```

A.1 The extended colour dataset

The parameters used in the Spriteworld software to generate the Shapes dataset with more (rainbow) colours:

```
COLOURS = {
  'red': distributions.Mixture([distributions.Continuous('c0', 0.95, 1.), distributions.Continuous('c0', 0., 0.05)]),
  'orange': distributions.Continuous('c0', 0.85, 0.10),
  'yellow': distributions.Continuous('c0', 0.10, 0.10),
  'green': distributions.Continuous('c0', 0.27, 0.37),
  'blue': distributions.Continuous('c0', 0.55, 0.65),
  'indigo': distributions.Continuous('c0', 0.68, 0.78),
  'violet': distributions.Continuous('c0', 0.75, 0.85),
}
```

B Neural architectures and hyper-parameters

image width	64
image height	64
image channels	3
CNN kernel size	4×4
CNN stride	2×2
CNN layers	4
CNN filters	64
CNN dense layers	2
CNN dense layer size	256
dimensions of latent space	6
initialization interval for means of priors	$[-1.0, 1.0]$
initialization interval for log-variances of priors	$[-7.0, 0.0]$
batch size	32
Adam learning rate	10^{-3}
Adam β_1	0.9
Adam β_2	0.999
Adam ϵ	10^{-7}

Author Contributions Sean Tull developed the mathematical formalisation and wrote the theory sections. Razin A. Shaikh wrote the code, ran the experiments, and prepared some of the figures. Sara Sabrina Zemljic created the data and helped run the experiments. Stephen Clark oversaw the project, ran some of the experiments, and wrote the remainder of the manuscript. All authors took part equally in setting the general direction of the project.

Funding N/A

Declarations

Competing interests The authors declare no competing interests.

Ethics approval N/A

Consent to participate N/A

Consent for publication Yes

References

- Abramson J, Ahuja A, Barr I, Brussee A, Carnevale F, Cassin M (2020) DeepMind-interactive-agents-group. Imitating Interactive Intelligence [arXiv:2012.05672](https://arxiv.org/abs/2012.05672)
- Aerts D (2009) Quantum structure in cognition. *J Math Psychol* 53(5):314–348
- Aerts D, Gabora L (2005) A state-context-property model of concepts and their combinations I: the structure of the sets of contexts and properties. *Kybernetes* 34:151–175
- Aisbett J, Gibbon G (2001) A general formulation of conceptual spaces as a meso level representation. *Artif Intell* 133(1–2):189–232
- Bechberger L, Kühnberger K-U (2017) A thorough formalization of conceptual spaces. Joint German/Austrian conference on artificial intelligence (künstliche intelligenz) (pp 58–71)
- Benedetti M, Lloyd E, Sack S, Fiorentini M (2019) Parameterized quantum circuits as machine learning models. *Quantum Sci Technol* 4(043001)
- Bengio Y, Courville A, Vincent P (2013) Representation learning: a review and new perspectives. *IEEE Trans Pattern Anal Mach Intell*
- Birkhoff G, von Neumann J (1936) The logic of quantum mechanics. *Ann Math* 37(4):823–843
- Bogacz R (2017) A tutorial on the free-energy framework for modelling perception and learning. *J Math Psychol* 76:198–211
- Bolt J, Coecke B, Genovesi F, Lewis M, Marsden D, Piedeleu R (2019) Interacting conceptual spaces I: grammatical composition of concepts. *Conceptual spaces: elaborations and applications*, Springer (pp 151–181)
- Bražinskas A, Havrylov S, Titov I (2018) Embedding words as distributions with a Bayesian skip-gram model. Proceedings of the 27th international conference on computational linguistics (pp 1775–1789). Santa Fe, New Mexico, USA: Association for Computational Linguistics. Retrieved from <https://aclanthology.org/C18-1151>
- Cho K, Jacobs B (2019) Disintegration and Bayesian inversion via string diagrams. *Math Struct Comput Sci* 29(7):938–971
- Cho K, Jacobs B, Westerbaan B, Westerbaan A (2015) An introduction to effectus theory. [arXiv:1512.05813](https://arxiv.org/abs/1512.05813)
- Clark S, Lerchner A, von Glehn T, Tieleman O, Tanburn R, Dashevskiy M, Bosnjak M (2021) Formalising concepts as grounded abstractions (Tech. Rep.). <https://arxiv.org/pdf/2101.05125.pdf>: DeepMind, London
- Coecke B (2006) Introducing categories to the practicing physicist. *What is Category Theory* 30:45–74
- Coecke B, Kissinger A (2017) *Picturing quantum processes: a first course in quantum theory and diagrammatic reasoning*. Cambridge University Press
- Coecke B, Sadzadeh M, Clark S (2010) Mathematical foundations for a compositional distributional model of meaning. [arXiv:1003.4394](https://arxiv.org/abs/1003.4394)
- Doersch C (2016) Tutorial on variational autoencoders (Tech. Rep.), UC Berkeley. [arXiv:1606.05908](https://arxiv.org/abs/1606.05908)
- Epping GP, Busemeyer JR (2022) Using diverging predictions from classical and quantum models to dissociate between categorization systems. <https://doi.org/10.31234/osf.io/fq2k5>
- Epping GP, Fisher EL, Zeleznikow-Johnston A, Pothos E, Tsuchiya N (2021) A quantum geometric framework for modeling color similarity judgements. <https://doi.org/10.31234/osf.io/vtzrq>
- Fong B (2019) *An invitation to applied category theory - seven sketches in compositionality*. Cambridge University Press

- Friston K, Kiebel S (2009) Predictive coding under the free-energy principle. *Philosophical transactions of the Royal Society B: Biological sciences* 364(1521):1211–1221
- Ganter B, Obiedkov S (2016) *Conceptual exploration*. Springer
- Ganter B, Wille R (1999) *Formal concept analysis: mathematical foundations*. Springer Science & Business Media
- Gärdenfors P (2004) *Conceptual spaces: the geometry of thought*. MIT press
- Gärdenfors P (2014) *The geometry of meaning*. The MIT Press
- Goodfellow I, Bengio Y, Courville A (2016) *Deep learning*. The MIT Press
- Goodwin GP, Johnson-Laird PN (2013) The acquisition of Boolean concepts. *Trends Cognit Sci* 17. <https://doi.org/10.1016/j.tics.2013.01.007>
- Gopnik A, Meltzoff A (1997) *Words, thoughts, and theories*. MIT Press
- Harnad S (1990) The symbol grounding problem. *Physica D: Nonlinear Phenomona* 42:335–346
- Havlicek V, Corcoles AD, Temme K, Harrow AW, Kandala A, Chow JM, Gambetta JM (2019) Supervised learning with quantum-enhanced feature spaces. *Nature* 567:209–212
- Higgins I, Matthey L, Pal A, Burgess CP, Glorot X, Botvinick M, Lerchner A (2017) β -VAE: learning basic visual concepts with a constrained variational framework. *Proceedings of ICLR 2017*
- Higgins I, Sonnerat N, Matthey L, Pal A, Burgess CP, Bošnjak M, Lerchner A (2018) SCAN: learning hierarchical compositional visual concepts. *Proceedings of ICLR 2018*
- Huang Q, Smolensky P, He X, Deng L, Wu D (2018) Tensor product generation networks for deep NLP modeling. *Proceedings of the 2018 conference of the north American chapter of the association for computational linguistics: Human language technologies, vol 1 (long papers)* (pp 1263–1273). New Orleans, Louisiana: Association for Computational Linguistics. Retrieved from <https://aclanthology.org/N18-1114>
- Khoshaman A, Vinci W, Denis B, Andriyash E, Sadeghi H, Amin MH (2018) Quantum variational autoencoder. *Quantum. Sci Technol* 4(1):014001
- Kingma DP, Welling M (2014) Auto-encoding variational Bayes. *Proceedings of the international conference on learning representations (ICLR 2014)*
- Lake BM, Ullman TD, Tenenbaum JB, Gershman SJ (2017) Building machines that learn and think like people. *Behav Brain Sci* 40
- Lewis M, Lawry J (2016) Hierarchical conceptual spaces for concept combination. *Artif Intell* 237:204–227
- Locatello F, Bauer S, Lucic M, Rätsch G, Gelly S, Schölkopf B, Bachem O (2019) Challenging common assumptions in the unsupervised learning of disentangled representations. *Proceedings of the 36th international conference on machine learning*. Long Beach, California
- Lorenz R, Pearson A, Meichanetzidis K, Kartsaklis D, Coecke B (2023) QNLP in practice: running compositional models of meaning on a quantum computer. *J Artif Intell Res* 76. <https://doi.org/10.1613/jair.1.14329>
- Margolis E, Laurence S (Eds.) (2015) *The conceptual mind: new directions in the study of concepts*. The MIT Press
- Margolis E, Laurence S (2022) *Concepts*. <https://plato.stanford.edu/archives/fall2022/entries/concepts/>. (The Stanford Encyclopedia of Philosophy)
- Murphy GL (2002) *The big book of concepts*. The MIT Press
- Panangaden P (1998) Probabilistic relations. *School of Computer Science Research Reports-University of Birmingham CSR* 59–74
- Pothos EM, Busemeyer JR (2013) Can quantum probability provide a new direction for cognitive modeling? *Behav Brain Sci* 36(3)
- Preskill J (2012) Quantum computing and the entanglement frontier. (Rapporteur talk at the 25th Solvay Conference on Physics - The Theory of the Quantum World). [arXiv:1203.5813](https://arxiv.org/abs/1203.5813)
- Rezende DJ, Mohamed S, Wierstra D (2014) Stochastic backpropagation and approximate inference in deep generative models. *Proceedings of the 31st international conference on machine learning* (pp 1278–1286)
- Rickard JT, Aisbett J, Gibbon G (2007) Reformulation of the theory of conceptual spaces. *Inf Sci* 177(21):4539–4565
- Rodatz B, Shaikh RA, Yeh L (2021) Conversational negation using worldly context in compositional distributional semantics. [arXiv:2105.05748](https://arxiv.org/abs/2105.05748)
- Rosch EH (1973) Natural categories. *Cognit Psychol* 4(3):328–350
- Schlangen D, Zarriß S, Kennington C (2016) Resolving references to objects in photographs using the words-as-classifiers model. *Proceedings of the 54th annual meeting of the association for computational linguistics (vol. 1: Long papers)* (pp 1213–1223). Berlin, Germany: Association for Computational Linguistics. Retrieved from <https://aclanthology.org/P16-1115>
- Schuld M, Killoran N (2019) Quantum machine learning in feature Hilbert spaces. *Phys Rev Lett* 122:040504. <https://doi.org/10.1103/PhysRevLett.122.040504>
- Selinger P (2010) A survey of graphical languages for monoidal categories. *New structures for physics*, Springer (pp 289–355)
- Shaikh RA, Yeh L, Rodatz B, Coecke B (2021) Composing conversational negation. [arXiv:2107.06820](https://arxiv.org/abs/2107.06820)
- Shiebler D, Gavranovic B, Wilson P (2021) Category theory in machine learning. *The 4th international conference on applied category theory*. Cambridge, UK
- Smolensky P, Legendre G (2006) *The harmonic mind*. The MIT Press
- Tomas V, Sylvie D (2015) Unitary transformations in the quantum model for conceptual conjunctions and its application to data representation. *Front Psychol* 6
- Trueblood JS, Busemeyer JR (2011) A quantum probability account of order effects in inference. *Cognit Sci* 35:1518–1552
- Tull S (2019) *Categorical operational physics*. [arXiv:1902.00343](https://arxiv.org/abs/1902.00343)
- Tull S (2021) A categorical semantics of fuzzy concepts in conceptual spaces. *Proceedings of Applied Category Theory 2021*
- Van de Wetering J (2021) Constructing quantum circuits with global gates. *New J Phys* 23(4):043015
- Watters N, Matthey L, Borgeaud S, Kabra R, Lerchner A (2019) Sprite-world: a flexible, configurable reinforcement learning environment. <https://github.com/deepmind/spriteworld/>. Retrieved from <https://github.com/deepmind/spriteworld/>
- Yan F, Li N, Hirota K (2021) Qhsl: a quantum hue, saturation, and lightness color model. *Inf Sci* 577:196–213

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.