

Article

Learning Numerosity Representations with Transformers: Number Generation Tasks and Out-of-Distribution Generalization

Tommaso Boccato ¹, Alberto Testolin ^{1,2,*}  and Marco Zorzi ^{1,3,*} 

¹ Department of General Psychology, University of Padova, Via Venezia 8, 35131 Padova, Italy; tommaso.boccatto@studenti.unipd.it

² Department of Information Engineering, University of Padova, Via Gradenigo 6, 35131 Padova, Italy

³ IRCCS San Camillo Hospital, Via Alberoni 70, 30126 Venice-Lido, Italy

* Correspondence: alberto.testolin@unipd.it (A.T.); marco.zorzi@unipd.it (M.Z.)

Abstract: One of the most rapidly advancing areas of deep learning research aims at creating models that learn to disentangle the latent factors of variation from a data distribution. However, modeling joint probability mass functions is usually prohibitive, which motivates the use of conditional models assuming that some information is given as input. In the domain of numerical cognition, deep learning architectures have successfully demonstrated that approximate numerosity representations can emerge in multi-layer networks that build latent representations of a set of images with a varying number of items. However, existing models have focused on tasks requiring to conditionally estimate numerosity information from a *given image*. Here, we focus on a set of much more challenging tasks, which require to conditionally generate synthetic images containing a *given number* of items. We show that attention-based architectures operating at the pixel level can learn to produce well-formed images approximately containing a specific number of items, even when the target numerosity was not present in the training distribution.

Keywords: deep neural networks; attention mechanisms; density estimation; numerosity perception; cognitive modeling



Citation: Boccato, T.; Testolin, A.; Zorzi, M. Learning Numerosity Representations with Transformers: Number Generation Tasks and Out-of-Distribution Generalization. *Entropy* **2021**, *23*, 857. <https://doi.org/10.3390/e23070857>

Academic Editor: Mohamed Medhat Gaber

Received: 17 May 2021
Accepted: 29 June 2021
Published: 3 July 2021

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

In recent years, there has been a growing interest in the challenging problem of unsupervised representation learning [1]. Compared to the first wave of supervised deep learning success [2], unsupervised learning has great potential to further improve the capability of artificial intelligence systems, since it would allow building high-level, flexible representations without the need of explicit human supervision. Unsupervised deep learning models are also plausible from a cognitive [3] and biological [4] perspective, because they suggest how the brain could extract multiple levels of representations from the sensory signal by learning a hierarchical generative model of the environment [5–8].

Early approaches based on deep belief networks [9] already established that unsupervised representation learning leads to the discovery of high-level visual features, such as object parts [10] or written shapes [11,12]. However, the full potential of deep generative models was revealed by the introduction of variational autoencoders (VAE) [13] and generative adversarial networks (GAN) [14], which can discover and factorize extremely abstract attributes from the data [15,16]. These architectures can be further extended to promote the emergence of even more disentangled representations, such as in beta-VAE [17] and InfoGAN [18], or can exploit attention mechanisms to produce meaningful decompositions of complex visual scenes [19].

An interesting case study to investigate the representational capability of deep learning models is that of *numerosity perception*, which consists of rapidly estimating the number

of objects in a visual scene without resorting to sequential counting procedures [20]. Compared to other high-level visual features, numerosity information is particularly challenging to extract because it refers to a global property of the visual scene, which co-varies with many other non-numerical visual features such as cumulative area, density and item size [21]. The emergence of numerosity representations has been successfully simulated using deep belief networks [22–24], which can approximately estimate the number of items in a given image (matching human-level performance) and partially disentangle it from non-numerical magnitudes [25]. However, learning fully disentangled representations of numerosity seems to be still out of reach even for state-of-the-art generative models, such as the InfoGAN [26].

In this paper, we investigate whether the deployment of *self-attention mechanisms* allows more precisely encoding numerosity information as a disentangled factor of variation. Attention mechanisms [27] were first introduced in the context of machine translation to overcome the limitations of sequence-to-sequence architectures [28], which aimed at compressing the information contained in temporal sequences into fixed-length latent vectors. Shortly after, a novel architecture based solely on attention called *Transformer* [29] achieved new heights by completely dropping recurrence and convolutions. In analogy with the dynamics of associative memories [30], the power of this approach lies in the possibility of using a global attention mechanism to precisely and adaptively weight the contribution of each input element during processing. Transformers are starting to also be applied outside the language domain, with notable success in challenging computer vision tasks [31–33].

These promising results motivated the present work, whose main goal is to demonstrate that attention mechanisms can be successfully exploited to learn disentangled representations of numerosity, which can be used to generate novel synthetic images approximately containing a given number of items. Inspired by recent approaches that evaluated the capability of deep generative models to create novel attributes and their combinations [34], the Transformer was probed in different generative scenarios requiring to produce specific numerosities that were never encountered during training. The internal structure of the representational code was also analyzed, in order to investigate whether numerosity information could be mapped into a lower dimensional space that preserves the semantics of cardinal numbers [35].

2. Methods

2.1. Problem Formulation

Let $\mathcal{D} = \{(x_1, n_1), \dots, (x_m, n_m)\}$ s.t. $(x, n) \sim p(x, n)$, $n \in \mathcal{N}$ be a training dataset consisting of images paired with their respective numerosity (i.e., the number of items contained in each image). The generic (x, n) tuple is sampled i.i.d. from the $p(x, n)$ joint probability mass function (PMF), with $\mathcal{N} \subset \mathbb{N}$. The goal is to model the $p(x|n)$ conditional PMF, exploiting a density estimation algorithm which relies solely on global attention applied to the raw input images. Modeling such density by disentangling numerosity from other factors of variation should ideally allow generating images with a controlled number of objects, which could be specified by manipulating the initial state of the generative process (although the generative model does not receive explicit knowledge about cardinal numbers during training, the initial state of the generative process is the same for all training images featuring the same numerosity, as explained in Section 2.2). Crucially, the generative model might even learn to produce out-of-distribution samples belonging to areas of the $p(x, n)$ support that are not represented in \mathcal{D} , that is, images containing a number of objects that was never experienced during training.

Practically, one could focus on an equivalent representation of the target density, which exploits the chain rule to allow the density estimation algorithm to work autoregressively:

$$p(x|n) = \prod_{i=1}^r p(x_i|x_1, \dots, x_{i-1}, n); \quad (1)$$

where $\mathbf{x} = (x_1, \dots, x_r)$ represents a flattened image x made by r pixels. Let $q(\mathbf{x}|n, \boldsymbol{\theta}^*)$ be the approximated conditional PMF produced by the density estimation algorithm, with $\boldsymbol{\theta}^*$ denoting the optimal model parameters; it originates from the minimization of the following negative log-likelihood:

$$L(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, n) \sim \mathcal{D}} [-\log q(\mathbf{x}|n, \boldsymbol{\theta})] \quad (2)$$

$$= \mathbb{E}_{(\mathbf{x}, n) \sim \mathcal{D}} \left[-\log \prod_{i=1}^r q(x_i|x_1, \dots, x_{i-1}, n, \boldsymbol{\theta}) \right] \quad (3)$$

$$= \mathbb{E}_{(\mathbf{x}, n) \sim \mathcal{D}} \left[-\sum_{i=1}^r \log q(x_i|x_1, \dots, x_{i-1}, n, \boldsymbol{\theta}) \right]. \quad (4)$$

Step (4) suggests that the Transformer can be straightforwardly trained by computing the CrossEntropyLoss criterion on the model output logits exploiting PyTorch [36].

2.2. Model Architecture

The model investigated in this work is an encoder-only Transformer capable of dealing with data characterized by spatial relationships (e.g., images); its backbone, indeed, is built from the self-attention layers devised in [29]. Overall, the following mapping is implemented: $(\mathbf{x}, n) \mapsto \mathbf{P} \in \mathbb{R}^{q \times p}$, where $\mathbf{x} = (x_1, \dots, x_q)$ denotes the categorical input intensities, $q \leq r$ and \mathbf{p}_i^T (i.e., the i th row of \mathbf{P}) represents the conditional PMF associated to x_i (the density support, of cardinality p , coincides with the set of input intensities).

Input frames are not directly fed into the Transformer encoder: pixel intensities are first scanned following the raster order, and then transformed into learnable embeddings to which the position information is added. Being the positional encodings also learned, it is important to highlight that the model is invariant with respect to the order in which inputs are supplied; however, once the order is fixed, it must be maintained. During the last processing stage, the encoder output goes through a linear layer. Hence, the conditional probability mass functions (PMFs) are computed by applying a softmax function to the produced logits.

The deployed encoder only accepts sequences of real-valued d -dimensional vectors. As a consequence, the supplied dataset entries undergo careful processing. Firstly, the (x_1, \dots, x_{q-1}) intensities are mapped into $q-1$ embeddings (pixel x_q is never consumed by autoregression), $\mathbf{X} \in \mathbb{R}^{(q-1) \times d}$. Then, the encoder input is computed as:

$$\mathbf{H}_0 = [\mathbf{s}, \mathbf{X}^T]^T + \mathbf{E}; \quad (5)$$

where $\mathbf{s} \in \mathbb{R}^d$ encodes the equivalence class to which the considered image belongs (i.e., the numerosity n) while $\mathbf{E} \in \mathbb{R}^{q \times d}$ stores information about the pixel positions. Borrowing the machine translation nomenclature, we call \mathbf{s} the *Start of String* (SoS). Input embeddings, SoSs and positional encodings are obtained in the same way: the discrete starting values (i.e., intensities, numerosities and positions) trivially become indexes capable of selecting the corresponding rows in one of the $\mathbf{W}_X \in \mathbb{R}^{p \times d}$, $\mathbf{W}_s \in \mathbb{R}^{|\mathcal{N}| \times d}$ and $\mathbf{W}_E \in \mathbb{R}^{r \times d}$ matrices, learned through backpropagation [31–33]. We emphasize that the learned embeddings minimize the introduction of explicit inductive biases: the untrained model, indeed, is completely unaware of the distances between gray levels, the ordering on \mathbb{N} and the correlations of pixels. As a side effect, \mathbf{W}_E constrains the input image resolution.

The encoder consists of $2L$ properly stacked multi-head scaled dot-product attention ($\text{mha}(\bullet)$) and point-wise fully connected ($\text{fc}(\bullet)$) sub-layers. Residual connections and layer normalizations ($\text{norm}(\bullet)$) complete the architecture. Resuming from (5),

$$\mathbf{A}_l = \text{norm}(\mathbf{H}_{l-1} + \text{mha}(\mathbf{H}_{l-1})) \quad (6)$$

$$\mathbf{H}_l = \text{norm}(\mathbf{A}_l + \text{fc}(\mathbf{A}_l)) \quad (7)$$

describe the encoder pipeline, with the $l \in [1, L]$ subscript denoting the considered layer. The detailed implementations of $\text{mha}(\bullet)$, $\text{fc}(\bullet)$ and $\text{norm}(\bullet)$ can be found in [29]. Finally, the $\text{linear}(\bullet)$ and $\text{softmax}(\bullet)$ functions are assembled to produce the target conditional densities:

$$P = \text{softmax}(\text{linear}(\mathbf{H}_L)). \quad (8)$$

The *attention graphs* [37] shown in Figure 1 help us in explaining how the encoder autoregression is achieved. Directed edges identify the allowed attention flows; the missing ones (with respect to the respective fully connected, bipartite sub-graphs) are masked to prevent queries from attending to illegal positions. Solid edges denote the active attention flows involved in the generation of the considered gray level. The generative loop starts as in the left graph, where the represented forward pass results in the sampling of the first intensity, \tilde{x}_1 . In the middle and right graphs, the gray level obtained during the previous pass is appended to the input sequence, and the process is repeated. Further details about the model architecture and training hyperparameters are reported in Appendix A.

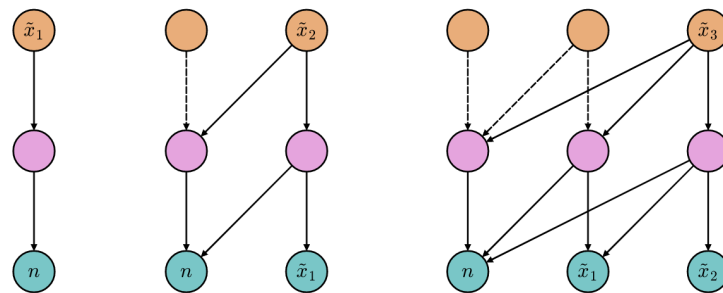


Figure 1. Example attention graphs ($L = 2$, $r = 3$) describing the *spontaneous generation* (see Section 2.4) of a novel image. Teal, pink and orange nodes represent the input, hidden and output positions, respectively. For each query node, the outgoing connections indicate which positions can be adaptively weighted.

2.3. Datasets

The Transformer was trained on two different datasets containing images of size 32×32 pixels with a varying number of objects (white dots) placed on a black background. Numerosities were uniformly sampled from the set $\{1, \dots, 8\}$. Each dataset was split into training, validation and test subsets containing, respectively, 16,000, 3200 and 3200 images. We verified that the size of the training set was properly calibrated by plotting the validation loss as a function of the number of training patterns (see Figure A2).

The first dataset, which we call *Uniform Dots*, contained images featuring objects of uniform size (see samples in the top row of Figure 2). In this dataset, the numerosity information is perfectly correlated with the total number of active pixels, which does not allow assessing to what extent the Transformer can disentangle numerosity from cumulative area. We thus also introduced a second dataset, which we call *Non-Uniform Dots*, containing images featuring objects of different size and constant (on average) cumulative area (see samples in bottom row of Figure 2). Let $A_{dot} \sim N(\mu_{dot}, \sigma_{dot}^2)$ be the random variable quantifying the individual area covered by a dot. The total area covered in a frame characterized by n dots can be expressed as:

$$A_{frame} = \sum_{i=1}^n A_{dot} \quad (9)$$

$$= nA_{dot} \quad (10)$$

$$\sim N(n\mu_{dot}, n^2\sigma_{dot}^2). \quad (11)$$

Setting $\mu_{dot} = \frac{\mu_{frame}}{n}$ implies that $A_{frame} \sim N(\mu_{frame}, n^2\sigma_{dot}^2)$, thus making the expected cumulative area $\mathbb{E}[A_{frame}] = \mu_{frame}$ independent from n (in our case, these parameters were set to $\mu_{frame} = 150$ and $\sigma_{dot} = 8$).

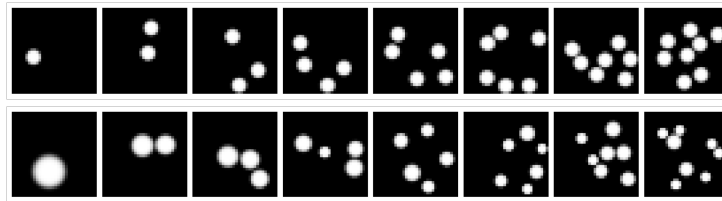


Figure 2. Sample images from the two datasets considered. Stimuli with increasing numerosity (i.e., $\mathcal{N} = \{1, \dots, 8\}$) are progressively shown from left to right. The top row contains samples from the *Uniform Dots* dataset, where the rendered dots have the same radius. The bottom row contains samples from the *Non-Uniform Dots*, where the area of each dot is sampled from $N(\mu_{frame}/n, \sigma_{dot}^2)$.

2.4. Generative Tasks

To investigate the emergence of numerosity representations, we designed a variety of generation tasks. The objective of these experiments was twofold: on the one hand, they allowed establishing whether the learned representations could be used to produce synthetic images with controlled properties (i.e., featuring a specific numerosity); on the other hand, they allowed studying the internal structure of the Transformer's latent space, in order to investigate whether it could embed the semantics of cardinal numbers.

As an initial assessment, the Transformer was evaluated in a straightforward *conditional generation* task: given the ground-truth (x, n) tuple, the goal is to approximate x through the modeled $q(x|n, \theta^*)$, incrementally building the image \tilde{x} according to:

$$\tilde{x}_i = \arg \max_x q(x|x_1, \dots, x_{i-1}, n, \theta^*), \quad \forall i \in [1, r]. \quad (12)$$

In other words, each pixel is determined by those preceding it in the fixed scan order, and the current ground-truth pixel values are provided as input at each time step. This task was only used to monitor learning progress, since it is well-known that one-step-ahead prediction is much easier than autoregressive self-generation [38].

In the more challenging *spontaneous generation* tasks, the Transformer was required to build an entire novel image \tilde{x} from scratch according to:

$$\tilde{x}_i \sim q(x|\tilde{x}_1, \dots, \tilde{x}_{i-1}, n, \theta^*), \quad \forall i \in [1, r]. \quad (13)$$

Unlike (12), each pixel intensity is now conditioned on the previously sampled ones; Equation (13), therefore, requires r forward passes for each image. It should be noted that, during the first generative step, the encoder input sequence contains only the SoS. Carrying on, the sequence gradually incorporates the new intensity embeddings: $s^T, [s, X_1^T]^T, \dots, [s, X_{r-1}^T]^T$, where the rows of X_i correspond to the first i sampled gray levels. For each SoS considered, a fixed number of 64 images was generated, in order to collect statistics about the samples produced. Spontaneous generation was tested under four different conditions:

- *Spontaneous generation over trained numerosities.* In this case, the sampling process was initially conditioned on the $|\mathcal{N}|$ numerosity representations learned during the Transformer training (i.e., the rows of W_s). That is, the learned SoSs were provided as initial seed.
- *Spontaneous generation over interpolated numerosities.* In this case, we tested whether the generative process could be biased toward specific numerosities that were never encountered during training (but nevertheless fell in the training interval) by injecting a novel SoS as initial seed. Defining w_i^T as the row of W_s corresponding to the training numerosity i , the desired conditioning n is injected by simply setting $s = (w_{n-1} + w_{n+1})/2$. In other words, the new representation of n is linearly interpolated from the two closest SoSs.

- *Spontaneous generation over extrapolated numerosities.* The generative capability was further pushed by exploring whether the Transformer could be biased to produce numerosities falling outside the training range. The proposed extrapolation mechanism relies on the *attribute vector* technique described in [39], where the vector representing the direction of change is computed as $\mathbf{a} = \mathbf{w}_{|\mathcal{N}|} - \mathbf{w}_{|\mathcal{N}|-1}$; it represents the direction along which the largest numerosities grow. We conjecture that the representation of a numerosity immediately larger than those included in the training range $[1, |\mathcal{N}|]$ can be approximated by $\mathbf{s} = \mathbf{w}_{|\mathcal{N}|} + \alpha \mathbf{a}$, for a suitable $\alpha > 0$.
- *Spontaneous generation with reduced components.* Although the embedding size is constrained by the encoder architecture, numerosity information might in fact be mapped into a lower-dimensional space, akin to an ordered “number line” [40]. To explore the possibility that the learned SoSs could be arranged along a one- or two-dimensional subspace, we performed a principal component analysis (PCA) on the rows of \mathbf{W}_s and used either the first or the first and second principal components to reconstruct the SoSs used to start the generation process and thus establish whether the sampling quality is affected by such dimensionality reduction.

After all image pixels are generated, the number of dots produced needs to be estimated using a suitable heuristic. For the purpose, two dataset-specific heuristics were introduced. The first counter is a simple area-based heuristic designed to work with uniform dots-like samples. The generated numerosity, indeed, can be computed by simply dividing the area covered by the rendered dots in a frame by the average dot area; such mean value is trivially estimated from the validation split of the dots dataset. Since this heuristic does not work in the case of dots with different size, to estimate the number of dots produced by the Transformer trained on the *Non-Uniform Dots* dataset, a ResNet18 classifier [41] was employed. The ResNet18 was trained on a subset (22,000 samples) characterized by $\mathcal{N} = \{0, \dots, 10\}$. Both counters achieved 100% accuracy on the respective dataset testing splits: the perfect accuracy achieved by the ResNet18 classifier on the test set suggests that numerosity estimation up to 10 items can be a trivial task for supervised deep learning models, at least with respect to the stimulus space considered in our simulations.

3. Results

After each generation task, the SoS-specific histograms of the generated numerosities were computed. We provide two different histogram visualizations: one depicts the relative frequency of each generated numerosity [34,42], while a 2D histogram is used to reproduce the visualization often used in human behavioral studies [43].

The generation histograms related to the *spontaneous generation over trained numerosities* task are shown in Figure 3. Especially for the *Uniform Dots* dataset (top panels), it is evident that the Transformer is able to create synthetic images with a specified numerosity, although the number of generated items is not always accurate. The generation is almost perfect for very small numbers (i.e., 1 and 2), while the model often generates one extra or one fewer item when asked to produce images with larger numerosities (see also the sample images reported in Figure A1). A similar pattern of errors is observed when the Transformer is trained using the *Non-Uniform Dots* dataset (bottom panels), although in this case the sampling uncertainty associated with larger numerosities increases, and the model sometimes generates images with a mismatch of up to three items. Overall, these results are well-aligned with the existing empirical literature on human behavior, which suggests that numerosity estimates are distributed around the target mean and variability tends to increase with numerosity [42,44], and that numerosity estimation can be altered by confounding non-numerical magnitudes [21,25].

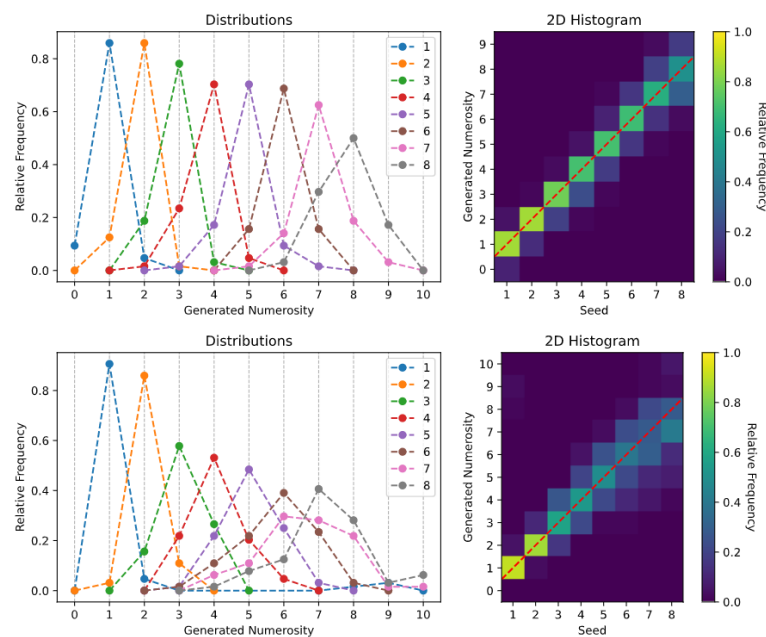


Figure 3. Spontaneous generation over trained numerosities: (Top) *Uniform Dots* dataset; and (Bottom) *Non-Uniform Dots* dataset.

Notably, the synthetic images produced by the Transformer are much more precise compared to samples produced by other deep generative models, such as VAEs or GANs [26,34]. Moreover, differently from previous approaches, here we demonstrate that the generation process can be biased toward a specific numerosity, suggesting that attention mechanisms play a key role in allowing a more precise processing of numerosity information. As a control simulation, we also tested the generative capability of a model trained on images containing objects of a different shape. To this aim, we created the *Smoothed Squares* dataset containing images produced by inscribing squares into the circles of the *Uniform Dots* dataset, applying an average filtering (3×3) and a gamma correction ($x_{out} = x_{in}^{0.25}$). Sample images from this dataset are shown in Figure A3. Histograms related to the spontaneous generation task are shown in Figure A4, while samples of generated images are reported in the left panel of Figure A5. Interestingly, the Transformer generates well-formed images even when trained on a dataset containing a mix of images from *Smoothed Squares* and *Uniform Dots* (see the right panel of Figure A5), suggesting that it can properly factorize also shape information.

The generation histograms related to the *spontaneous generation over interpolated numerosities* task are shown in the left panel of Figure 4. Quite impressively, the Transformer is able to produce images with a specific number of objects even for numerosities that were never encountered during training. For example, by averaging the embeddings corresponding to $n = 1$ and $n = 3$, the model always generates images with exactly two dots (orange line in the left panel). An analogous finding holds when interpolating the numerosities $n = 4$ and $n = 6$, although in those cases the number of items is not always perfectly matched (sample images are reported in Figure A6). These remarkable findings suggest that the emergent representational space approximately encodes the semantics of cardinal numbers, at least within the lower and upper training bounds.

As shown in the right panel of Figure 4, the results related to the *spontaneous generation over extrapolated numerosities* task further corroborate this hypothesis. Indeed, the attribute vector computed as the difference between the embeddings of the two largest numerosities in the training set (in this case, $n = 4$ and $n = 5$) seems to represent the direction of increase of the numerosity feature: by summing a fraction of such vector to the embedding of $n = 5$, the Transformer can reliably generate images containing six items, although sometimes the additional item appears squeezed or slightly distorted (sample images are reported

in Figure A7). Interestingly, when the attribute vector is scaled by a factor $\alpha = 0.5$, the Transformer equally generates images with either five or six items. However, setting $\alpha \geq 2$ did not allow to reliably generate images with seven items, suggesting that the learned embeddings approximately capture a sort of “successor function” only over the local neighborhood of a specific numerosity.

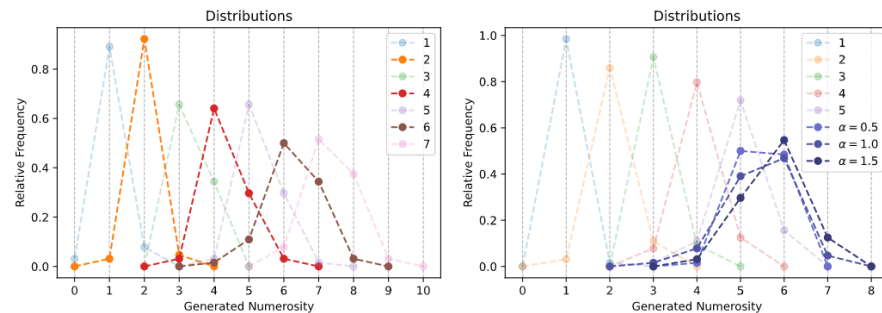


Figure 4. Spontaneous generation over interpolated (left) and extrapolated (right) numerosities. Trained numerosities are represented by semi-transparent curves, while solid curves represent unseen numerosities.

The lower-dimensional manifold structure of the encoder space is shown in Figure 5. Interestingly, and in partial alignment with other recent computational work [35], it seems that the topology of the numerosity embeddings preserves the strict ordering of cardinal numbers, even though the Transformer did not explicitly receive such information during training. This is evident even by just looking at the first principal component (x-axis in the figure), which suggests that numerosity information could be internally organized as a one-dimensional “number line” [20,45]. However, differently from Kondapaneni and Perona [35], we found that the second principal component does not monotonically encode cardinal information, but suggests a periodic pattern. As a control analysis, Figure A8 also shows the PCA projection resulting right after the random initialization of the embeddings, which indeed does not reflect any ordering structure.

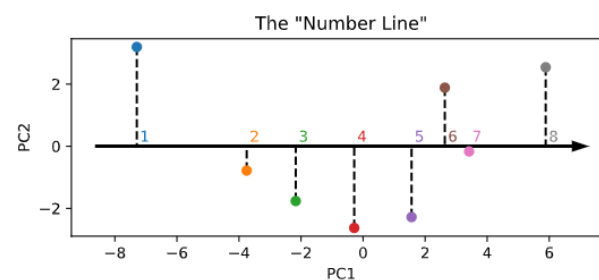


Figure 5. Visualization of the lower-dimensional embedding space resulting from PCA.

The results related to the *spontaneous generation with reduced components* task suggest that, when the embeddings are projected into such lower-dimensional manifold, the generative abilities of the model are preserved: as shown in the top panels of Figure 6, the Transformer can generate samples with remarkable accuracy even when only the first principal component is retained. Adding the second principal component (bottom panels of Figure 6) allows further improving the generation precision, although numerosities mapped to nearby points in the lower dimensional space (i.e., $n = 6$ and $n = 7$) are frequently confounded.

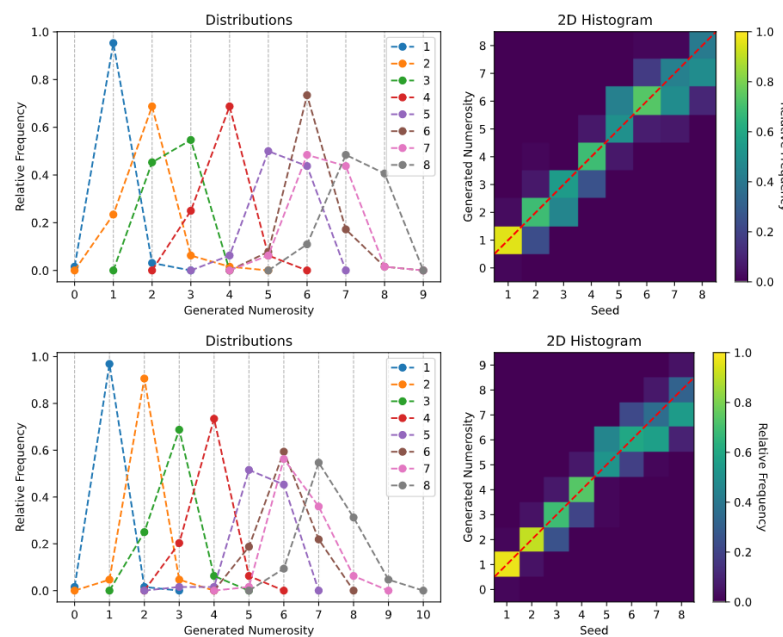


Figure 6. Spontaneous generation with reduced components, considering: one principal component (**top**); and two principal components (**bottom**).

4. Conclusions

In this study, we investigated whether state-of-the-art deep learning architectures based on attention mechanisms could learn disentangled representations of numerosity from a set of images containing a variable number of items. Our simulations not only show that Transformers can successfully learn to generate synthetic images featuring a target numerosity, but also that they can interpolate and extrapolate the generation process to previously unseen numerosities. These remarkable findings suggest that Transformers can indeed disentangle numerosity from other non-numerical visual features. However, it should be stressed that the generation process is error-prone and thus reflects an *approximate* representation of numerical information. Moreover, although we are impressed by the Transformer’s generative capabilities, in real world scenarios, the number of training patterns can be exponentially smaller than the support of the probability mass function to be estimated, which makes generalization to out-of-distribution samples particularly challenging [34]. A key open issue is thus to establish whether domain-general deep learning architectures could extrapolate numerical knowledge well beyond the limit of their training distribution, which would require learning more abstract conceptual structures, such as the successor function [46], which form the foundation of our understanding of natural numbers [47].

Another limitation of Transformer architectures is related to their computational complexity: naive implementations have a quadratic cost in the number of pixels in terms of both memory and computation, preventing their scaling to high resolutions. Recent studies have attempted to mitigate this issue by approximating global attention in different ways, for example by restricting self-attention receptive fields to local neighborhoods [31], reducing image resolution [32] or focusing on image patches [33]. In the present work, the image size allowed efficiently training and testing the Transformer architecture; however, future work should better clarify whether more effective attention mechanisms could be employed to scale-up the model to realistic image sizes.

Author Contributions: Conceptualization and supervision, A.T. and M.Z.; methodology, T.B. and A.T.; software, simulations, analysis, and visualization, T.B.; writing—original draft preparation, T.B. and A.T.; writing—review and editing, T.B., A.T. and M.Z.; project administration, M.Z.; and funding acquisition, A.T. and M.Z. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by Cariparo Foundation (Excellence Grant 2017 “NUMSENSE” to M.Z.)

Data Availability Statement: The source code for the simulations is available for download at <https://github.com/BoCtrl-C?tab=repositories> (accessed on 30 June 2021).

Conflicts of Interest: The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

Appendix A. Model Hyperparameters and Supplementary Figures

The Transformer described in the present work is available in three different sizes: “S” (~37,000 parameters), “M” (~136,000 parameters) and “L” (~568,000 parameters); Table A1 reports the corresponding hyperparameters. For publication purposes, all presented results refer to the size “M” model. However, we also investigated the performance of the small and large variants on a subset of the introduced tasks, without noticing major differences. All models were trained using Adam optimizer [48] ($\beta_1 = 0.9$, $\beta_2 = 0.999$) with a batch size of 16. The training always stops when the validation loss does not improve, with respect to the best loss obtained, for five epochs in a row. After a search in the set $\{0.03, 0.01, 0.003, 0.001\}$, the initial learning rate was set to 0.003. Furthermore, the learning rate decays by 0.1 every 25 training epochs. Each training session took, on average, 1 h 45 min on a workstation equipped with an NVidia GTX 1080 graphic card. To reduce the Transformer computational demand, the dataset entries were pre-processed by uniformly quantizing (16 levels) the input intensities (i.e., $p = 16$).

Table A1. Hyperparameters characterizing the available model sizes. See the PyTorch Transformer documentation for more details about `nhead` and `dim_feedforward` (<https://pytorch.org/docs/stable/generated/torch.nn.Transformer.html> (accessed on 30 June 2021)).

Hyperparameter	Value		
	Size “S”	Size “M”	Size “L”
<code>d_model</code> (d)	16	32	64
<code>nhead</code>	2	2	4
<code>num_encoder_layers</code> (L)	6	8	10
<code>dim_feedforward</code>	64	128	256

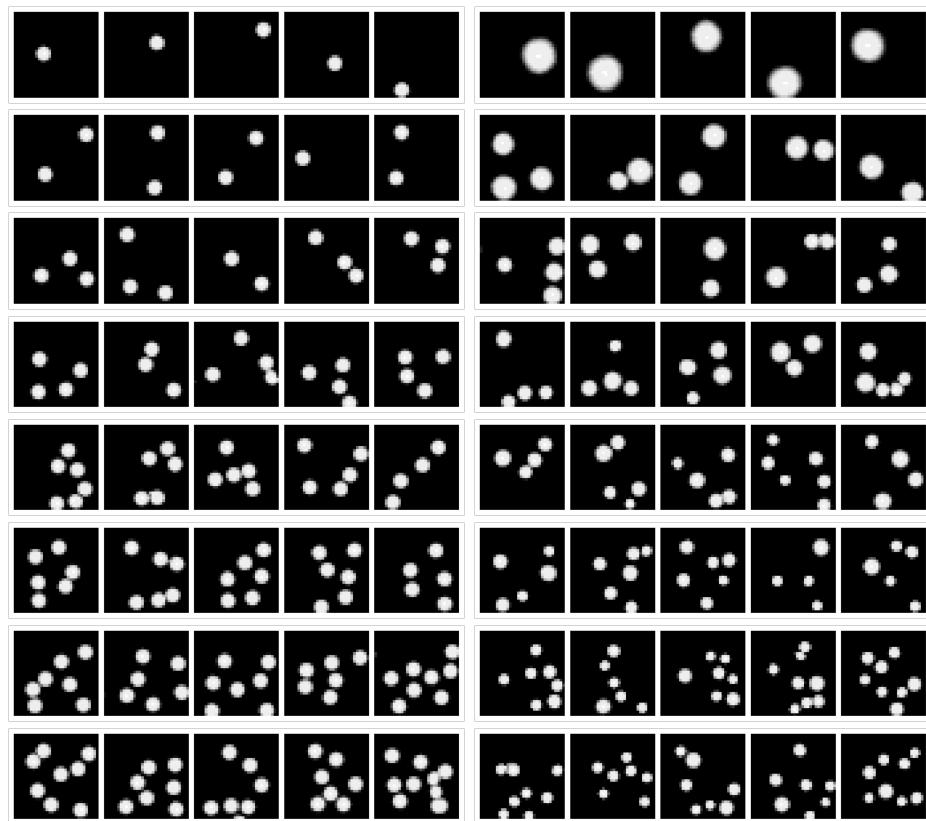


Figure A1. Sample images produced in the *spontaneous generation over trained numerosities* task: **(Left)** *Uniform Dots*; and **(Right)** *Non-Uniform Dots*. Images produced with increasing generation seeds (i.e., $\mathcal{N} = \{1, \dots, 8\}$) are progressively shown from top to bottom. Sample images with a number of dots that does not match the seed are purposely included for illustration.

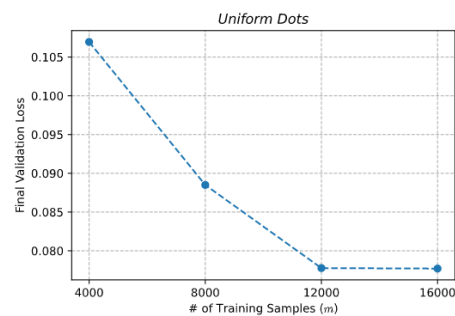


Figure A2. Converging validation loss obtained by training the Transformer (size “M”) on instances of the *Uniform Dots* dataset characterized by an increasing number of training patterns.

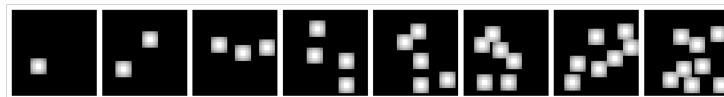


Figure A3. Sample images from the *Smoothed Squares* dataset. Stimuli with increasing numerosity (i.e., $\mathcal{N} = \{1, \dots, 8\}$) are progressively shown from left to right.

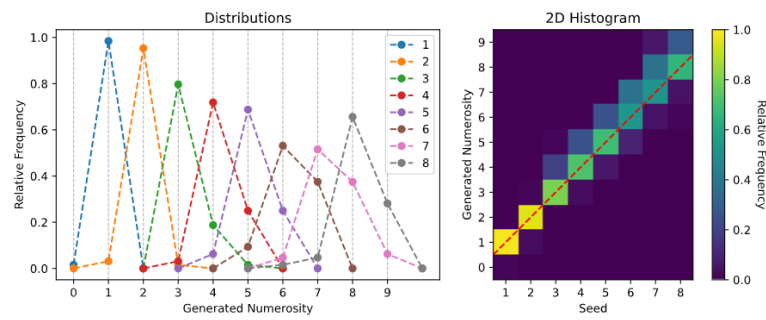


Figure A4. Spontaneous generation over trained numerosities on the *Smoothed Squares* dataset.

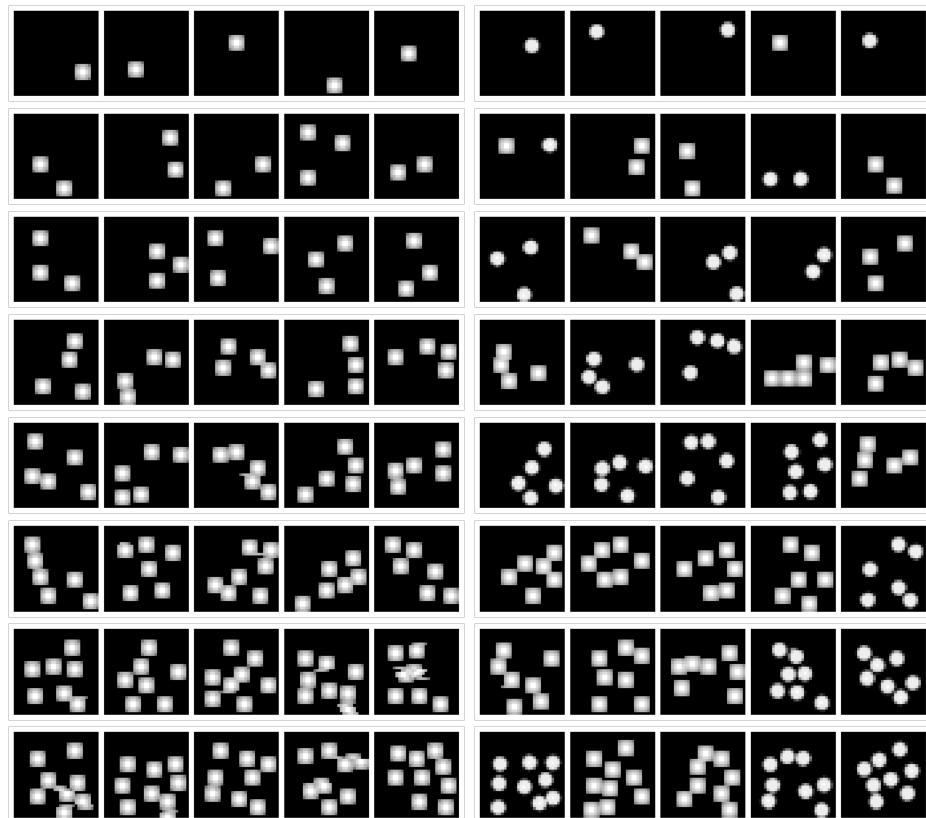


Figure A5. Sample images produced in the *spontaneous generation over trained numerosities* task by the model trained on *Smoothed Squares* (left) or a mix of *Smoothed Squares* and *Uniform dots* (right). Images produced with increasing generation seeds (i.e., $\mathcal{N} = \{1, \dots, 8\}$) are progressively shown from top to bottom.

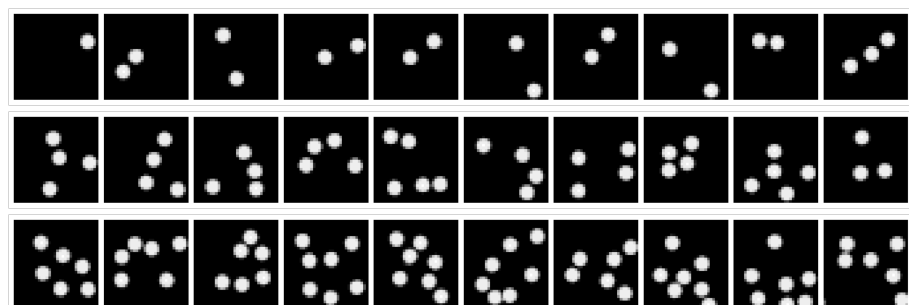


Figure A6. Sample images conditioned on the interpolated SoSs (i.e., *spontaneous generation over interpolated numerosities* task). From top to bottom, the displayed rows correspond to the unseen numerosities 2, 4 and 6, respectively.

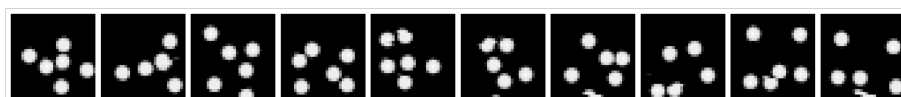


Figure A7. Sample images conditioned on the extrapolated SoS (i.e., *spontaneous generation over extrapolated numerosities* task). The results reported refer to $\alpha = 1$, which should produce images containing six objects. Note how the rendering of dots is qualitatively less precise than the ones shown in Figures A1 and A6.

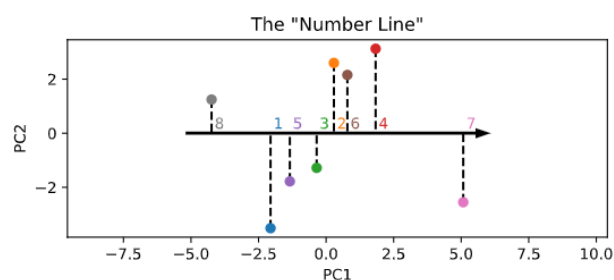


Figure A8. Visualization of the lower-dimensional embedding space right after random initialization.

References

- Bengio, Y.; Courville, A.; Vincent, P. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.* **2013**, *35*, 1798–1828. [[CrossRef](#)] [[PubMed](#)]
- LeCun, Y.; Bengio, Y.; Hinton, G. Deep learning. *Nature* **2015**, *521*, 436–444. [[CrossRef](#)] [[PubMed](#)]
- Zorzi, M.; Testolin, A.; Stoianov, I.P. Modeling language and cognition with deep unsupervised learning: a tutorial overview. *Front. Psychol.* **2013**, *4*, 515. [[CrossRef](#)] [[PubMed](#)]
- Zhuang, C.; Yan, S.; Nayebi, A.; Schrimpf, M.; Frank, M.C.; DiCarlo, J.J.; Yamins, D.L. Unsupervised neural network models of the ventral visual stream. *Proc. Natl. Acad. Sci. USA* **2021**, *118*, e2014196118. [[CrossRef](#)]
- Clark, A. Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behav. Brain Sci.* **2013**, *36*, 181–204. [[CrossRef](#)]
- Friston, K. The free-energy principle: A unified brain theory? *Nat. Rev. Neurosci.* **2010**, *11*, 127–138. [[CrossRef](#)]
- Hinton, G.E. Learning multiple layers of representation. *Trends Cogn. Sci.* **2007**, *11*, 428–434. [[CrossRef](#)]
- Testolin, A.; Zorzi, M. Probabilistic models and generative neural networks: Towards a unified framework for modeling normal and impaired neurocognitive functions. *Front. Comput. Neurosci.* **2016**, *10*, 73. [[CrossRef](#)]
- Hinton, G.E.; Osindero, S.; Teh, Y.W. A fast learning algorithm for deep belief nets. *Neural Comput.* **2006**, *18*, 1527–1554. [[CrossRef](#)]
- Le, Q.V. Building high-level features using large scale unsupervised learning. In Proceedings of the 2013 IEEE international conference on acoustics, speech and signal processing, Vancouver, BC, Canada, 26–31 May 2013; pp. 8595–8598.
- Sadeghi, Z.; Testolin, A. Learning representation hierarchies by sharing visual features: a computational investigation of Persian character recognition with unsupervised deep learning. *Cogn. Process.* **2017**, *18*, 273–284. [[CrossRef](#)]
- Testolin, A.; Stoianov, I.; Zorzi, M. Letter perception emerges from unsupervised deep learning and recycling of natural image features. *Nat. Hum. Behav.* **2017**, *1*, 657–664. [[CrossRef](#)]
- Kingma, D.P.; Welling, M. Auto-encoding variational bayes. *arXiv* **2013**, arXiv:1312.6114.
- Goodfellow, I.J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; Bengio, Y. Generative adversarial networks. *arXiv* **2014**, arXiv:1406.2661.
- Karras, T.; Laine, S.; Aila, T. A style-based generator architecture for generative adversarial networks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, Long Beach, CA, USA, 16–20 June 2019; pp. 4401–4410.
- Zhu, J.Y.; Park, T.; Isola, P.; Efros, A.A. Unpaired image-to-image translation using cycle-consistent adversarial networks. In Proceedings of the IEEE International Conference on Computer Vision, Venice, Italy, 22–29 October 2017; pp. 2223–2232.
- Higgins, I.; Matthey, L.; Pal, A.; Burgess, C.; Glorot, X.; Botvinick, M.; Mohamed, S.; Lerchner, A. beta-vae: Learning Basic Visual Concepts with a Constrained Variational Framework 2016. Available online: <https://openreview.net/forum?id=Sy2fzU9gl> (accessed on 15 May 2021)
- Chen, X.; Duan, Y.; Houthoofd, R.; Schulman, J.; Sutskever, I.; Abbeel, P. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *arXiv* **2016**, arXiv:1606.03657.
- Burgess, C.P.; Matthey, L.; Watters, N.; Kabra, R.; Higgins, I.; Botvinick, M.; Lerchner, A. Monet: Unsupervised scene decomposition and representation. *arXiv* **2019**, arXiv:1901.11390.
- Dehaene, S. *The Number Sense: How the Mind Creates Mathematics*; Oxford University Press: Oxford, England, 2011.

21. Gebuis, T.; Reynvoet, B. The interplay between nonsymbolic number and its continuous visual properties. *J. Exp. Psychol. Gen.* **2012**, *141*, 642. [[CrossRef](#)] [[PubMed](#)]
22. Stoianov, I.; Zorzi, M. Emergence of a 'visual number sense' in hierarchical generative models. *Nat. Neurosci.* **2012**, *15*, 194–196. [[CrossRef](#)]
23. Testolin, A.; Zou, W.Y.; McClelland, J.L. Numerosity discrimination in deep neural networks: Initial competence, developmental refinement and experience statistics. *Dev. Sci.* **2020**, *23*, e12940. [[CrossRef](#)]
24. Zorzi, M.; Testolin, A. An emergentist perspective on the origin of number sense. *Philos. Trans. R. Soc. B Biol. Sci.* **2018**, *373*, 20170043. [[CrossRef](#)]
25. Testolin, A.; Dolfi, S.; Rochus, M.; Zorzi, M. Visual sense of number vs. sense of magnitude in humans and machines. *Sci. Rep.* **2020**, *10*, 1–13.
26. Zanetti, A.; Testolin, A.; Zorzi, M.; Wawrzynski, P. Numerosity Representation in InfoGAN: An Empirical Study. In Proceedings of the International Work-Conference on Artificial Neural Networks, Gran Canaria, Spain 12–14 June 2019; pp. 49–60.
27. Bahdanau, D.; Cho, K.; Bengio, Y. Neural Machine Translation by Jointly Learning to Align and Translate. In Proceedings of the 3rd International Conference on Learning Representations, San Diego, CA, USA, 7–9 May 2015.
28. Sutskever, I.; Vinyals, O.; Le, Q.V. Sequence to Sequence Learning with Neural Networks. In *Advances in Neural Information Processing Systems*; Ghahramani, Z.; Welling, M., Cortes, C., Lawrence, N., Weinberger, K.Q., Eds.; Curran Associates Inc.: Red Hook, NY, USA, 2014; Volume 27.
29. Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, L.; Polosukhin, I. Attention is All you Need. In *Advances in Neural Information Processing Systems*; Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R., Eds.; Curran Associates, Inc.: Red Hook, NY, USA, 2017; Volume 30.
30. Ramsauer, H.; Schäfl, B.; Lehner, J.; Seidl, P.; Widrich, M.; Gruber, L.; Holzleitner, M.; Pavlović, M.; Sandve, G.K.; Greiff, V.; et al. Hopfield networks is all you need. *arXiv* **2020**, arXiv:2008.02217.
31. Parmar, N.; Vaswani, A.; Uszkoreit, J.; Kaiser, L.; Shazeer, N.; Ku, A.; Tran, D. Image Transformer. In Proceedings of the 35th International Conference on Machine Learning, Stockholm, Sweden, 10–15 July 2018; Volume 80, pp. 4055–4064.
32. Chen, M.; Radford, A.; Child, R.; Wu, J.; Jun, H.; Luan, D.; Sutskever, I. Generative Pretraining From Pixels. In Proceedings of the 37th International Conference on Machine Learning, Virtual Event, 13–18 July 2020; Volume 119, pp. 1691–1703.
33. Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *arXiv* **2020**, arXiv:2010.11929.
34. Zhao, S.; Ren, H.; Yuan, A.; Song, J.; Goodman, N.; Ermon, S. Bias and Generalization in Deep Generative Models: An Empirical Study. In Proceedings of the Advances in Neural Information Processing Systems, Montréal, QC, Canada, 3–8 December 2018; Volume 31.
35. Kondapaneni, N.; Perona, P. A Number Sense as an Emergent Property of the Manipulating Brain. *arXiv* **2020**, arXiv:2012.04132.
36. Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems, Vancouver, BC, Canada, 8–14 December 2019*; Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., Garnett, R., Eds.; Curran Associates, Inc.: Red Hook, NY, USA, 2019; Volume 32.
37. Abnar, S.; Zuidema, W. Quantifying Attention Flow in Transformers. *arXiv* **2020**, arXiv:cs.LG/2005.00928.
38. Cenzato, A.; Testolin, A.; Zorzi, M. Long-Term Prediction of Physical Interactions: A Challenge for Deep Generative Models. In *International Conference on Machine Learning, Optimization, and Data Science Siena, Italy, 10–13 September 2019*; Springer: Berlin/Heidelberg, Germany, 2020; pp. 83–94.
39. Carter, S.; Nielsen, M. Using Artificial Intelligence to Augment Human Intelligence. *Distill* **2017**. [[CrossRef](#)]
40. Dehaene, S. The neural basis of the Weber–Fechner law: A logarithmic mental number line. *Trends Cogn. Sci.* **2003**, *7*, 145–147. [[CrossRef](#)]
41. He, K.; Zhang, X.; Ren, S.; Sun, J. Deep Residual Learning for Image Recognition. In Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, NV, USA, 27–30 June 2016; pp. 770–778. [[CrossRef](#)]
42. Sella, F.; Berteletti, I.; Lucangeli, D.; Zorzi, M. Spontaneous non-verbal counting in toddlers. *Dev. Sci.* **2016**, *19*, 329–337. [[CrossRef](#)] [[PubMed](#)]
43. Revkin, S.K.; Piazza, M.; Izard, V.; Cohen, L.; Dehaene, S. Does subitizing reflect numerical estimation? *Psychol. Sci.* **2008**, *19*, 607–614. [[CrossRef](#)]
44. Testolin, A.; McClelland, J.L. Do estimates of numerosity really adhere to Weber's law? A reexamination of two case studies. *Psychon. Bull. Rev.* **2021**, *28*, 158–168. [[CrossRef](#)]
45. Harvey, B.M.; Klein, B.P.; Petridou, N.; Dumoulin, S.O. Topographic representation of numerosity in the human parietal cortex. *Science* **2013**, *341*, 1123–1126. [[CrossRef](#)]
46. Leslie, A.M.; Gelman, R.; Gallistel, C. The generative basis of natural number concepts. *Trends Cogn. Sci.* **2008**, *12*, 213–218. [[CrossRef](#)] [[PubMed](#)]
47. Testolin, A. The challenge of modeling the acquisition of mathematical concepts. *Front. Hum. Neurosci.* **2020**, *14*. [[CrossRef](#)] [[PubMed](#)]
48. Kingma, D.P.; Ba, J. Adam: A Method for Stochastic Optimization. In Proceedings of the 3rd International Conference on Learning Representations (ICLR), San Diego, CA, USA, 7–9 May 2015.