
Visual Enumeration is Challenging for Large-scale Generative AI

Alberto Testolin

Department of General Psychology and
Department of Mathematics
University of Padova, Padova (IT)
alberto.testolin@unipd.it

Kuinan Hou

Department of General Psychology
University of Padova, Padova (IT)
kuinan.hou@phd.unipd.it

Marco Zorzi

Department of General Psychology and Padova Neuroscience Center
University of Padova, Padova (IT) and
IRCSS San Camillo Hospital, Venice-Lido (IT)
marco.zorzi@unipd.it

ABSTRACT

Humans can readily judge the number of objects in a visual scene, even without counting, and such a skill has been documented in many animal species and babies prior to language development and formal schooling. Numerical judgments are error-free for small sets, while for larger collections responses become approximate, with variability increasing proportionally to the target number. This response pattern is observed for items of all kinds, despite variation in object features (such as color or shape), suggesting that our visual number sense relies on abstract representations of numerosity. Here, we investigate whether large-scale generative Artificial Intelligence (AI) systems have a human-like number sense, which should allow them to reliably name the number of objects in simple visual stimuli or generate images containing a target number of items in the 1-10 range. Surprisingly, most of the foundation models considered have a poor number sense: They make striking errors even with small numbers, the response variability does not increase in a systematic way, and the pattern of errors depends on object category. Only the most recent proprietary systems exhibit signatures of a visual number sense. Our findings demonstrate that having an intuitive visual understanding of number remains challenging for foundation models, which in turn might be detrimental to the perceptual grounding of numeracy that in humans is crucial for mathematical learning.

Keywords Foundation Models · Machine Vision · Deep Learning · BLIP · ViLT · GPT · Gemini · DALL-E · Stable Diffusion · Numerical Cognition · Numerosity · Counting

Introduction

Artificial Intelligence (AI) is progressing rapidly, with deep learning models approaching or even surpassing human performance in a variety of domains, including perceptual judgements [1] and natural language processing [2]. In this article, we investigate whether advanced AI systems can judge the numerosity of visual sets, a core capability that humans share with many animal species [3]. Even infants are sensitive to numerosity [4] and toddlers can generate sets containing a target number of items [5], suggesting that a preverbal understanding of numerical quantities develops well before formal education. Small numerosities in the “subitizing” range (up to 4) are perceived in an exact manner, while the numerosity of larger sets is approximately estimated when counting is precluded

[6]. In the latter case, responses follow Weber’s law, so that variability increases proportionally to the mean estimate [3]. Another key signature of number sense is its abstract nature: numerosity is encoded independently from object category, location or presentation modality [4]. Importantly, numerosity is spontaneously extracted by our visual system [7] and there is broad consensus that numerosity perception is foundational for subsequent learning of symbolic numbers as well as for the acquisition of higher-level mathematical competence [8, 3].

Machine vision researchers have engineered a variety of specialized systems for visual object counting, often tailored to specific categories such as penguins [9] or crowds [10]. The most popular framework consists on first running an object detector to segment the target items in a visual scene and then explicitly counting the resulting bounding boxes or object proposals [11, 12], often summing fractional counts estimated from different sections of the image [13]. However, in these approaches numerosity representations do not emerge within the model itself, because the encoding of number is delegated to an external, often hard-wired mechanism. A different perspective considers the possibility that numerosity representations might spontaneously emerge in deep neural networks as a high-order statistical feature of the sensory environment [14]. Indeed, it has been shown that a rudimentary visual number sense can emerge in small-scale generative models trained with the goal of reconstructing images with a varying number of objects [15, 16].

In this work, we investigate whether numerosity perception abilities spontaneously emerge in state-of-the-art generative AI systems. To this end, we systematically probe the visual number sense of several “foundation models”, which are large-scale generative architectures trained on huge data sets that have shown emergent abilities in a variety of domains [17] and can readily solve a wide range of downstream tasks [18]. Unlike the domain-specific architectures mentioned above, foundation models are domain-general systems that can be used out-of-the-box without the need of fine-tuning on numerical tasks. However, despite their flexibility and their remarkable performance in a variety of problems, it has been repeatedly shown that foundation models often fall short in tasks that require the manipulation of numerical information [19].

We investigate key properties of visual number sense across a range of models of different sizes and complexities. In the image-to-text domain, we consider modern Visual Question Answering (VQA) systems that can provide written answers to non-trivial questions about the content of an image or accurate descriptions of complex visual scenes. In particular, we test the capability of judging visual numerosity of the Vision-and-Language Transformer (ViLT) [20] and the Bootstrapping Language-Image Pre-training (BLIP-2) model [21]. We also consider two recent proprietary models, GPT-4V [22] developed by OpenAI and Gemini [23] developed by Google, which are considered the most advanced multimodal AI systems currently available. In the text-to-image domain, we consider generative models that can produce high-quality visual content following detailed user prompts provided in natural language. We test the numerosity production skills of popular generative architectures for images: Stable Diffusion (version 2.1) [24] and DALL-E (version 2 and version 3) [25, 26].

In line with the proposal of using methods from cognitive science to test foundation models [27], we exploit two behavioral tasks that are widely used to evaluate number sense in humans: numerosity naming [6], which requires establishing how many items are present in a given image, and numerosity production [28, 5], which requires generating a target number of items. We characterize the distribution of responses produced by the different AI models in the range 1-10 using a variety of object categories. Perfect accuracy across the entire numerical range would suggest the emergence of counting skills, while error-free responses with only small numbers would either indicate subitizing capabilities, or that counting is only partially developed as in children who do not fully master the counting principles [29, 30]. Error-prone responses centered on the target number would instead suggest that the AI model relies on approximate estimation mechanisms, which may or may not follow Weber’s law.

Results

Numerosity naming

We probed the numerosity naming skills of the Vision-and-Language Transformer (ViLT) [20], the Bootstrapping Language-Image Pre-training (BLIP-2) model [21], the latest multimodal version of the Generative Pre-trained Transformer (GPT-4V) [22] and the recently introduced multimodal Gemini model [23]. These systems have remarkable visual reasoning abilities and can answer non-trivial questions related to image content (e.g., *What does the image represent? What are the feelings of*

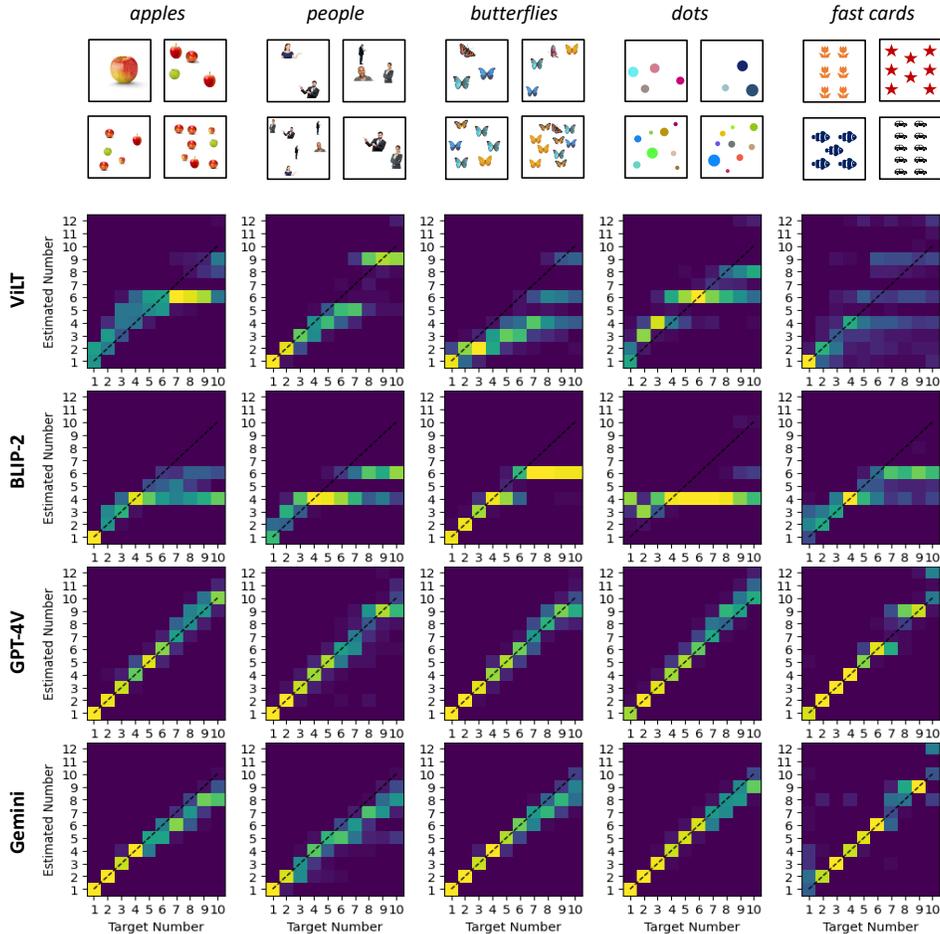


Figure 1: Confusion matrices for the numerosity naming task. Each panel shows the distribution of models’ responses across different object categories: apples, people, butterflies, dots and fast cards (a few sample stimuli are shown below each category name). The x-axis represents the target number, while the y-axis represents the corresponding model responses. Response frequency is encoded using a perceptually uniform colormap (blue = 0%, yellow = 100%).

the people in the scene and why?). We asked the models how many objects were present in a set of images containing up to 10 objects (see Materials and Methods for methodological details). Each stimulus included only items of the same kind, sampled from 5 possible classes: realistic pictures of common categories (apples, people, butterflies), colored dots, and “fast cards” depicting regularly placed clip-arts similar to those used to test children [29]. Examples of visual stimuli are shown at the top of Fig. 1.

In general, the response accuracy was extremely low for both ViLT and BLIP models (ViLT: 28.0%, BLIP-2: 29.6%), indicating that they cannot count. Moreover, in sharp contrast with humans, ViLT and BLIP models also returned wrong answers within the subitizing range (1-4) and even for images with one or two objects. According to standard criteria used in human developmental studies (see Methods), these models could be considered at most “One”-knowers, that is, they exhibit reliable enumeration only for a single object, as typical of children younger than three years of age [29]. Confusion matrices (CMs) reported in Fig. 1 clearly show the presence of anchoring effects, leading the models to choose stereotyped responses (e.g., 4, 6) probably corresponding to common spatial layouts that might be over-represented in their training corpora. The pattern of responses also drastically varied between categories (minimum correlation between CMs for ViLT: 0.04, BLIP-2: 0.27), suggesting that they fail to abstract numerical information.

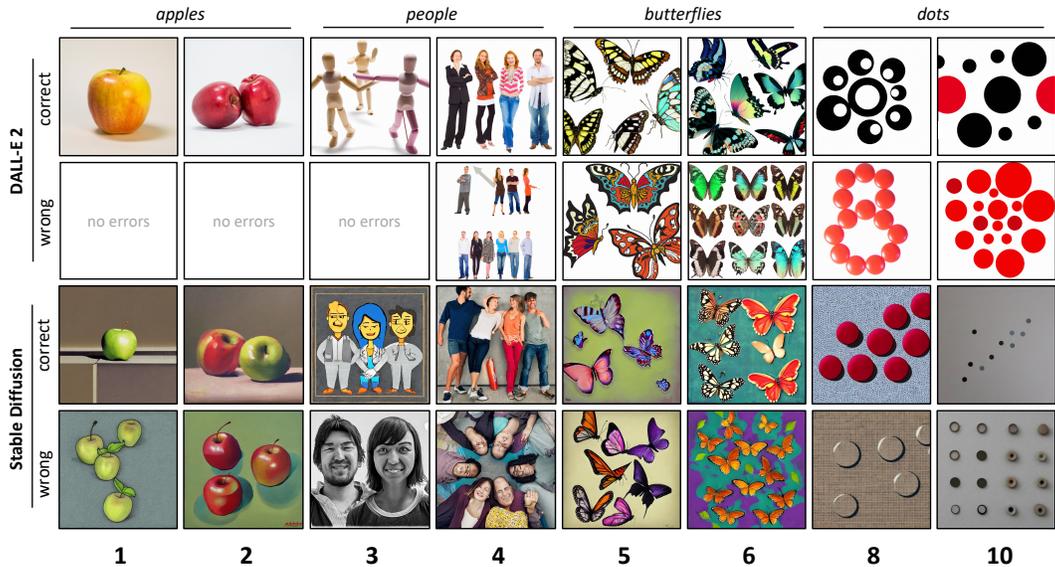


Figure 2: Examples of images generated by DALL-E 2 and Stable Diffusion in the numerosity production task, showcasing both correct and wrong generations (the target number is indicated at the bottom). We report two images for each target category: apples, people, butterflies, and dots. For some categories, DALL-E 2 never made mistakes when there were a few objects to produce. For the dots category, in a few cases DALL-E 2 generated images containing a wrong number of dots, which were nevertheless arranged according to the target digit shape (e.g., 8 in the figure).

The responses of GPT-4V and Gemini were much more accurate (73.9% and 70.6%, respectively), suggesting that these multimodal models might possess rudimentary enumeration skills. Responses were often error-free in the subitizing range, and confusion matrices were fairly consistent between categories (the minimum correlation was 0.87 for GPT-4V and 0.70 for Gemini). Considering the distribution of response errors, both GPT-4V and Gemini would be characterized as “Six”-knowers. Interestingly, such a pattern is never observed in human numeracy development because children typically transition from a “Four”-knower to a “Cardinal-principle”-knower level, which implies full mastery of the counting principles with no strict upper limit [29].

Numerosity production

We probed the numerosity production skills of three popular generative AI systems for images: Stable Diffusion (version 2.1) [24] and DALL-E (both versions 2 and 3) [25, 26]. These foundation models have proven capable of generating high-quality images following a textual description, also taking into account stylistic instructions, fine-grained details, and relational features (e.g., *A photo of an astronaut riding a horse in photorealistic style*). We asked the models to generate images with a target number of objects, in analogy with numerosity production tasks used in animal and human studies [28, 5]. Target objects belonged to the same classes used for the naming task, except for the “fast cards” category, which could be underrepresented in the corpora used to train these models. Examples of generated images are shown in Fig. 2, while confusion matrices are shown in Fig. 3.

Overall, DALL-E 2 tended to generate images with a white background and more definite objects, while Stable Diffusion adopted more artistic generation styles, producing either paintings, clip-arts, or realistic content. DALL-E 3 produced more detailed images with respect to its predecessor model. The mean response accuracy was fairly low for all models (Stable Diff: 33.3%, DALL-E 2: 38.7%, DALL-E 3: 47.7%). Interestingly, DALL-E 2 was the only model capable of exhibiting error-free responses in the small number range, but only in specific cases (up to 3 objects for the “apples” category, exactly 3 objects for the “people” category, and exactly 1 object for the “butterflies” category). Despite its higher average accuracy, DALL-E 3 always produced a few errors even in the subitizing range. According to criteria used in human developmental studies, Stable Diffusion would be considered a “Two”-knower, DALL-E 2 a “Three”-knower, and DALL-E 3 a “Four”-knower. Compared to the naming task, the

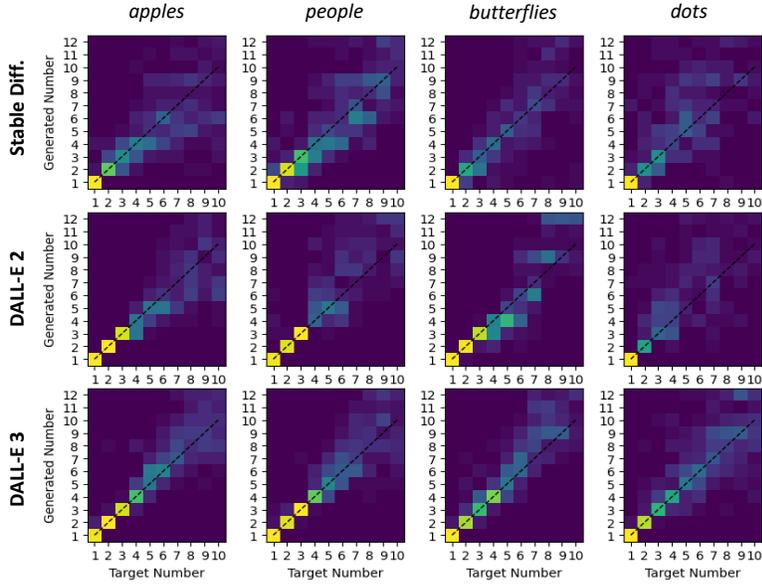


Figure 3: Confusion matrices for the numerosity production task. The x-axis represents the target number, while the y-axis represents the corresponding model responses. Response frequency is encoded using a perceptually uniform colormap (blue = 0%, yellow = 100%).

response patterns were more homogeneous across categories (minimum correlation between CMs for Stable Diff: 0.81, DALL-E 2: 0.76, DALL-E 3: 0.93).

Analysis of response distribution

In humans and other animal species, numerosity tasks yield a distribution of error responses that is not random, but varies systematically according to Weber’s law [31]. In particular, the standard error of the response increases proportionally to the mean, indicating scalar variability, which can be indexed by a constant coefficient of variation (standard error of the responses/mean response) [32, 33]. To investigate whether models’ responses followed a scalar variability pattern, we used a well-established methodology implementing a log–log regression analysis [34]. If scalar variability holds, as the number n increases the logs of the standard error $SE(n)$ and of the mean estimates $\mu(n)$ should form a straight line, with a slope $\beta = 1$ [$\log(SE(n)) = \alpha + \beta \log(\mu(n))$]. A slope that is significantly different from 1 would indicate that the relation between the standard errors and the means is a more complex power function, while a poor fit would more generally indicate that the response variability cannot be explained by linear or power trends.

We carried out a separate analysis for each model and for each object category. If a model exhibited error-free responses for small numbers, we excluded those numbers from the analysis. In many cases the regression analysis resulted in a poor fit ($p > 0.05$), indicating that the distribution of responses across numerosities did not match the trend observed in human and animal experiments. As detailed below, in the remaining cases the regression analysis was significant (all $p < 0.05$; $r^2 = 0.64 - 0.96$), so we further tested whether the regression slopes were significantly different from 1. Results were mixed, and depended both on the model type and the object category.

In the numerosity naming task, for ViLT and BLIP the distribution of models’ responses did not generally follow the pattern observed in human studies. For ViLT, we observed scalar variability only for the people (slope = 1.11, $p > 0.05$) and butterflies (slope = 1.14, $p > 0.05$) categories, while we observed power variability for the fast cards category (slope = 1.37, $p = 0.003$). For BLIP2, the regression analysis never produced a good fit. For GPT-4V, we observed scalar variability for the people (slope = 0.98, $p > 0.05$), butterflies (slope = 1.07, $p > 0.05$) and fast cards (slope = 1.23, $p > 0.05$) categories, while we observed power variability for the dots category (slope = 0.73, $p = 0.008$). For Gemini, the coefficient of variation was not stable across numerosities, resulting in a poor fit for almost

all categories ($p > 0.05$). Only for the people category we observed a systematic power variability trend (slope = 0.62, $p = 0.033$).

In the numerosity production task the errors were more regularly distributed around the target number. For Stable Diffusion, we observed scalar variability for the apples (slope = 1.07, $p > 0.05$) and butterflies (slope = 1.15, $p > 0.05$) categories, while we observed power variability for the people (slope = 0.70, $p = 0.030$) and dots (slope = 0.60, $p = 0.016$) categories. For DALL-E 2, we observed scalar variability for the dots category (slope = 0.82, $p > 0.05$), while we observed power variability for the apples (slope = 1.53, $p = 0.024$) and butterflies (slope = 1.49, $p = 0.036$) categories. DALL-E 3 was the model most closely adhering to Weber’s law, since we observed scalar variability across all categories: apples (slope = 0.98, $p > 0.05$), people (slope = 1.10, $p > 0.05$), butterflies (slope = 0.85, $p > 0.05$) and dots (slope = 0.98, $p > 0.05$).

Overall, scalar variability was thus mostly observed only in the most recent OpenAI models, GPT-4V and DALL-E 3. In the other models the response variability often increased according to a power law or did not have a systematic trend.

Discussion

This work demonstrates that large-scale foundation models cannot yet reliably enumerate the number of objects in a visual scene, both in image-to-text and text-to-image tasks. Such a striking deficit is observed even for sets containing just a few items, suggesting a number knower-level that is at best comparable to that of preschool children that do not fully master the counting principles [29, 30]. This might explain anecdotal evidence showing the failure of generative AI to synthesize realistic images featuring multiple instances of a body part, such as the correct number of fingers in a hand.

Surprisingly, even approximate number estimation failed to match the psychophysics of human numerosity perception. Only the most recent proprietary models, GPT-4V and DALL-E 3 from OpenAI, exhibited sparks of human-like number sense: their responses were often error-free for small numbers, suggesting some subitizing capabilities, and sometimes followed scalar variability for larger numbers, in accordance to Weber’s law. Also Gemini, the most recent multimodal model from Google, achieved a fairly accurate performance in the numerosity naming task, with error-free responses in the subitizing range. However, the response variability for larger numbers in this case did not adhere to Weber’s law.

The fact that the concept of visual numerosity remains elusive for large-scale foundation models is particularly striking when considering that numerosity in humans is susceptible to adaptation effects [35], which is the hallmark of a primary perceptual property (just as orientation or color). Accordingly, neuronal populations encoding numerosity have been found in multiple cortical regions in human neuroimaging experiments [36, 37, 38] and in neurophysiological studies with animals [39]. Moreover, computational modeling studies have shown that sensitivity to numerosity can emerge in small-scale deep learning models trained to generate synthetic images of object sets [15, 16, 40]. Diffusion models, such as Stable Diffusion and DALL-E, are trained with a similar objective on huge and heterogeneous image datasets that most likely contain substantial variability in numerosity. However, the empirical distribution of visual numerosities in natural image datasets is approximately captured by Zipf’s law (that is, frequency is inversely proportional to number [41]), therefore, it might be possible that oversampling of small numerosities in the training datasets of foundation models has detrimental effects on their emergent representational space. This still does not explain the variability across object categories observed in our analyses, which suggests that the representation of numerosity is not fully disentangled from other image properties. An alternative explanation for the poor enumeration skills of foundation models lies in the mapping between perceptual numerosity representations, encoded in image embeddings, and number symbols (number words or Arabic digits) encoded in text embeddings. In children, establishing such a bidirectional mapping is a sophisticated developmental process, which takes many years and requires explicit instruction [42].

In this respect, it would be valuable to analyze and compare the embeddings emerging in different models to gain a deeper understanding of the computational principles that could enable the development of stronger enumeration skills. Unfortunately, the closed-source nature of proprietary systems does not allow us to draw strong conclusions about the nature of their putative visual number sense. The fact that performance was never perfect even for the most advanced models suggests that numerosity estimation was not engineered in the system, but we cannot exclude the possibility that some counting mechanisms were partially introduced as extra processing layers during prompt elaboration. Furthermore, the most

sophisticated AI systems can also exploit the self-generation of code snippets to fulfill a user request, as in the case of mathematical problem solving [43]. In other words, an advanced AI system could in principle use external tools (e.g., based on object detection and a symbolic counting algorithm) to carry out visual enumeration without any understanding of numerosity and counting principles.

Whether numerosity representations could spontaneously develop as an emergent ability [17] in foundation models thus remains an important open question, which would require “opening the box” and inspecting the model’s inner functioning in a way that is not currently possible with proprietary models. This issue highlights the dangers of using proprietary models in academic research [44], and calls for further research efforts to develop open source foundation models with a focus on basic perceptual abilities, such as those underlying our visual number sense.

In conclusion, we believe that substantial progress in architecture design and training procedures is warranted to create AI systems that truly understand visual numerosity [45]. Grounding numeracy development in the visual number sense, as in the case of humans [3], might also be the key to enabling AI systems to acquire and fully master numerical and mathematical knowledge without resorting to highly specialized hybrid architectures [46].

Materials and Methods

In the naming task, for each object class and target number we probed the models using 50 high-resolution (1024×1024 pixels) images created by randomly placing items of variable size on a uniform white background, with no overlap. Fast card stimuli were created using clip-arts of common objects (apples, bells, butterflies, candies, cars, fish, flowers, planes, stars) drawn in different colors (black, blue, green, orange, red). Responses without numerical content or with vague quantification (e.g., “a few”) were automatically discarded, but we ensured that at least 20 acceptable answers were produced for each class/number combination. In the production task, for each condition we probed the models to generate 50 high-resolution images, which were then manually labeled by one of the authors (K.H.) and individually double checked by another author. Images with an ambiguous number of objects or with ill-formed content were discarded (see Materials and Methods for details). We ensured that at least 20 acceptable images were produced for each test condition.

Models

For the numerosity naming task, we considered two different Visual Question Answering (VQA) models. We used the vilt-b32-mlm version of the Vision-and-Language Transformer (ViLT) model [20] available through Hugging Face. This architecture incorporates text embeddings into a Vision Transformer (ViT), allowing it to have a minimal design for Vision-and-Language Pre-training and thus speeding-up model training and inference phases. It has a total of 87.4 million parameters. We also used the blip2-flan-t5-xl version of the Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models (BLIP-2) model [21], developed by Salesforce and available through Hugging Face. This architecture is considered an improved version of BLIP [47] and it is available in different versions, depending on the type and size of the backbone models. We explored all versions except for the t5-xxl model (due to GPU memory constraints) and found that the version using Flan-T5 as a language model yielded the best accuracy. The chosen model version has a total of 4.1 billion parameters.

We also considered the two most advanced multimodal models currently available, GPT-4V [22] and Gemini Pro [23], which can be readily used in VQA settings. These models are regarded among the most powerful generalist AI systems to date, thanks to their unprecedented ability to understand and process an arbitrary mix of input images and texts. Technical details regarding the underlying architecture and inner working of these models (including engineered modules that might be used to solve specific tasks) have not been revealed; it has been speculated that both these models might have more than one trillion of parameters¹.

For the numerosity production task, we considered two different image generation architectures. One is represented by the Stable Diffusion (version 2.1) model [24], developed by Stability AI and available through Hugging Face. Stable Diffusion v2.1 is a text-to-image latent diffusion model that combines an autoencoder with a diffusion model that is trained in the latent space of the autoencoder. It has

¹<https://the-decoder.com/gpt-4-architecture-datasets-costs-and-more-leaked>

approximately 500 million parameters. We also tested models from the DALL-E family (version 2 and version 3) [25, 26] using the API interface provided by OpenAI. DALL-E 2 is an improved version of the original text-to-image DALL-E model, featuring a total of 3.5 billion parameters. DALL-E 3 is the latest version, which was trained using highly descriptive, synthetic captions for the training images. The number of model parameters is currently unknown.

Prompting methods

For each model, we explored different prompting methods and selected the one leading to the best performance, measured as mean absolute distance from the target number.

For the numerosity naming task, we implemented three different prompting methods. The first required to explicitly estimate the number of objects belonging to a specific category (i.e., *How many apples / butterflies / people / dots / shapes are there in the picture?*). The other two used the more general “objects” or “things” words to identify the items to count. For ViLT, GPT-4V and Gemini, the best performance was achieved with the generic “things” prompt, while for BLIP-2 the best performance was achieved with the category-specific prompts.

To make sure that the models were prompted correctly, we carried out a control simulation related to a non-numerical task using the entire set of “apples” stimuli. We probed the VQA models with the following prompt: *What does the image represent?* and we considered as correct the following answers: *apple(s)* and *fruit*. All models almost always provided the correct answer across the entire set of stimuli (accuracy for ViLT : 96.8%, BLIP-2 : 100%, GPT-4V : 100%, Gemini: 100%), thus demonstrating a proper understanding of the image content and the prompt structure.

For the numerosity production task, both DALL-E and Stable Diffusion were initially prompted with the following text: *An image with n apples / butterflies / people / dots* (where n varied between 2 and 10). When $n = 1$ the prompt was adjusted to the singular form. However, for the dots category this prompting method resulted in poor generations: we obtained better results when the models were prompted with a more specific description of the image: *An image with white background with n filled dots*.

Analysis of model responses

For each model, we assessed the knower-level by applying standard criteria used in the literature on the development of counting skills [29] to the average responses across categories. To be considered an “ n ”-knower (i.e., “One”-knower, “Two”-knower, “Three”-knower, “Four”-knower) the model had to: 1) Give n objects at least 67% of the time when asked for that number; and 2) Give n objects no more than half as often when asked for a different number.

Responses given by VQA models were automatically parsed. If present, number words were converted to numerical values using the word2number Python library. The response was discarded if it contained multiple numbers or if it contained vague quantification terms (e.g., “a few”, “a bunch of”). Interestingly, in a few cases the Gemini model produced unexpected responses with images containing only one item, for example answering that “There are two things in the image: an apple and a white background”. We discarded these answers. We checked that at least 20 correct trials were recorded for each number / category. All models achieved the minimum number of tests required without the need to further prompt them (total number of responses discarded for ViLT: 0; BLIP-2: 6; GPT-4V: 0; Gemini: 20).

Images produced by the image generation models were discarded when the annotators judged them to be too ambiguous to be correctly parsed. The criteria for discarding an image were the following (representative examples are shown in Supplementary Figure S11):

- the model generated only objects of non-target categories (Fig. S11a);
- the model generated target objects that could in principle be identified as such, but were not well-formed because the shape of the item was significantly distorted and/or other distinguishing features were significantly altered (see examples *b* and *c* in Fig. S11);
- the model generated one or more objects that resembled the target category but could not be reliably identified as such in isolation (i.e., without providing the entire content of the image as context); this happened, for example, when an object was only partially visible (e.g., less than 10% of the object was included in the frame) and the model did not depict it using

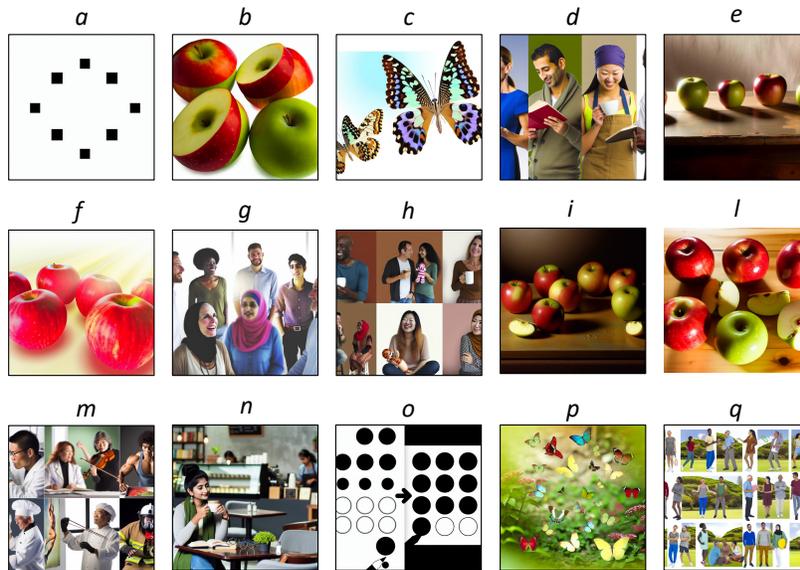


Figure S11: Examples of images produced by DALL-E 3 that illustrate the criteria we used to discard ambiguous generations. We also report the target number requested in the prompt to highlight the fact that sometimes the discarded image contained the correct amount of items, but often it did not. a) Wrong object category (target category: dots; target number = 6). b) Altered features or distorted object representations for the apples category (target number = 4). c) Altered features or distorted object representations for the butterflies category (target number = 3). d) Unreliable identification of a target object due to the absence of representative features (target number = 4). e) Unreliable identification of an object due to the absence of representative features (target number = 6). f) Reliable identification even though objects are only partially visible (target number = 6; counted 6 objects). g) Reliable identification even though objects are only partially visible (target number = 7; counted 7 objects). h) Reliable identification even though objects are only partially visible (target number = 7; counted 8 objects). i) Apples cut into pieces (target number = 6). l) Apples cut into pieces (target number = 8). m) Pictorial representations of people mixed with real people (target number = 6). n) People in the background (target number = 1). o) More than 20 objects (target number = 8). p) More than 20 objects (target number = 9). q) More than 20 objects (target number = 10).

representative features, such as the face of a person or the stalk of an apple (see examples *d* and *e* in Fig. S11); vice versa, objects were considered valid and thus counted when they appeared behind other objects or partially outside the image frame, but could nevertheless be reliably identified (see examples *f*, *g* and *h* in Fig. S11);

- for the "apple" category, we initially considered valid images containing unconventional representations of the target items (e.g., an apple cut in half); however, DALL-E 3 frequently generated images with entire apples plus several slices of apples and/or additional apples cut in half; to avoid ambiguity we thus opted for discarding those images (see examples *i* and *l* in Fig. S11);
- for the "people" category, we discarded images mixing realistic people and pictorial content representing people (e.g., an artist painting a portrait, see example *m* in Fig. S11), as well as images depicting people in the forefront but also less-visible people in the background (see example *n* in Fig. S11);
- we discarded images containing more than 20 objects (see examples *o*, *p* and *q* in Fig. S11).

It is important to clarify that each discarded image was replaced by an additionally generated sample. Moreover, the procedure was blind with respect to the target number of objects. This implies that the performance of a model cannot be penalized by the data cleaning procedure.

For Stable Diffusion, the minimum number was not achieved for the dots category with $n = 1$ (15), $n = 8$ (14), $n = 9$ (14) and $n = 10$ (13), since the model in these cases most frequently generated more than

30 items. For DALL-E 2, the minimum number of trials was not achieved for the dots category with $n = 10$ (15). Also in this case, for these prompts the model most frequently generated more than 30 items.

Software and computing hardware

All simulations and analyses were implemented using Python v3 and Google Colab. To test larger-scale models (BLIP-2 family) we used a virtual machine with an NVIDIA L4 GPU and 64 GB of RAM memory, which was allocated using the Google cloud computing platform. Data and code used in the current study are available on GitHub.

Acknowledgements

We are grateful to OpenAI for granting research access to the GPT-4V and DALL-E APIs. This work was partially supported by the Italian Ministry of University and Research (PRIN grant n. C53D23004110006 to M.Z.).

References

- [1] Scott Mayer McKinney, Marcin Sieniek, Varun Godbole, Jonathan Godwin, Natasha Antropova, Hutan Ashrafian, Trevor Back, Mary Chesus, Greg S Corrado, Ara Darzi, et al. International evaluation of an AI system for breast cancer screening. *Nature*, 577(7788):89–94, 2020.
- [2] Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. ChatGPT outperforms crowd-workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30), 2023.
- [3] Stanislas Dehaene. *The number sense: How the mind creates mathematics*. OUP USA, 2011.
- [4] Véronique Izard, Coralie Sann, Elizabeth S Spelke, and Arlette Streri. Newborn infants perceive abstract numbers. *Proceedings of the National Academy of Sciences*, 106(25):10382–10385, 2009.
- [5] Francesco Sella, Ilaria Berteletti, Daniela Lucangeli, and Marco Zorzi. Spontaneous non-verbal counting in toddlers. *Developmental science*, 19(2):329–337, 2016.
- [6] Susannah K Revkin, Manuela Piazza, Véronique Izard, Laurent Cohen, and Stanislas Dehaene. Does subitizing reflect numerical estimation? *Psychological science*, 19(6):607–614, 2008.
- [7] Guido Marco Cicchini, Giovanni Anobile, and David C Burr. Spontaneous perception of numerosity in humans. *Nature communications*, 7(1):12536, 2016.
- [8] Justin Halberda, Michèle MM Mazzocco, and Lisa Feigenson. Individual differences in non-verbal number acuity correlate with maths achievement. *Nature*, 455(7213):665–668, 2008.
- [9] Carlos Arteta, Victor Lempitsky, and Andrew Zisserman. Counting in the wild. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14*, pages 483–498. Springer, 2016.
- [10] Muhammad Asif Khan, Hamid Menouar, and Ridha Hamila. Revisiting crowd counting: State-of-the-art, trends, and future perspectives. *Image and Vision Computing*, 129:104597, 2023.
- [11] Alexander Trott, Caiming Xiong, and Richard Socher. Interpretable counting for visual question answering. In *International Conference on Learning Representations*, 2018.
- [12] Yan Zhang, Jonathon Hare, and Adam Prügel-Bennett. Learning to count objects in natural images for visual question answering. In *International Conference on Learning Representations*, 2018.
- [13] Prithvijit Chattopadhyay, Ramakrishna Vedantam, Ramprasaath R Selvaraju, Dhruv Batra, and Devi Parikh. Counting everyday objects in everyday scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1135–1144, 2017.
- [14] Marco Zorzi and Alberto Testolin. An emergentist perspective on the origin of number sense. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 373(1740):20170043, 2018.
- [15] Ivilin Stoianov and Marco Zorzi. Emergence of a ‘visual number sense’ in hierarchical generative models. *Nature neuroscience*, 15(2):194–196, 2012.

- [16] Alberto Testolin, Serena Dolfi, Mathijs Rochus, and Marco Zorzi. Visual sense of number vs. sense of magnitude in humans and machines. *Scientific reports*, 10(1):10045, 2020.
- [17] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022.
- [18] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [19] Alberto Testolin. Can neural networks do arithmetic? a survey on the elementary numerical skills of state-of-the-art deep learning models. *arXiv preprint arXiv:2303.07735*, 2023.
- [20] Wonjae Kim, Bokyung Son, and Ildoo Kim. ViLT: Vision-and-language transformer without convolution or region supervision. In *International Conference on Machine Learning*, pages 5583–5594, 2021.
- [21] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International Conference on Machine Learning*, 2023.
- [22] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The dawn of Imms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9(1):1, 2023.
- [23] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [24] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF conference on computer vision and pattern recognition*, 2022.
- [25] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- [26] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- [27] Marcel Binz and Eric Schulz. Using cognitive psychology to understand GPT-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120, 2023.
- [28] John Whalen, Charles R Gallistel, and Rochel Gelman. Nonverbal counting in humans: The psychophysics of number representation. *Psychological science*, 10(2):130–137, 1999.
- [29] Mathieu Le Corre and Susan Carey. One, two, three, four, nothing more: An investigation of the conceptual sources of the verbal counting principles. *Cognition*, 105(2):395–438, 2007.
- [30] Michael D Lee and Barbara W Sarnecka. Number-knower levels in young children: Insights from bayesian modeling. *Cognition*, 120(3):391–402, 2011.
- [31] Roger N Shepard, Dan W Kilpatrick, and James P Cunningham. The internal representation of numbers. *Cognitive psychology*, 7(1):82–138, 1975.
- [32] Stanislas Dehaene. The neural basis of the weber–fechner law: a logarithmic mental number line. *Trends in cognitive sciences*, 7(4):145–147, 2003.
- [33] Charles R Gallistel and Rochel Gelman. Non-verbal numerical cognition: From reals to integers. *Trends in cognitive sciences*, 4(2):59–65, 2000.
- [34] Véronique Izard and Stanislas Dehaene. Calibrating the mental number line. *Cognition*, 106(3):1221–1247, 2008.
- [35] David Burr and John Ross. A visual sense of number. *Current biology*, 18(6):425–428, 2008.
- [36] Ben M Harvey, Barrie P Klein, Natalia Petridou, and Serge O Dumoulin. Topographic representation of numerosity in the human parietal cortex. *Science*, 341(6150):1123–1126, 2013.

- [37] Elisa Castaldi, Manuela Piazza, Stanislas Dehaene, Alexandre Vignaud, and Evelyn Eger. Attentional amplification of neural codes for number independent of other quantities along the dorsal visual stream. *Elife*, 8:e45160, 2019.
- [38] Jacob M Paul, Martijn van Ackooij, Tuomas C Ten Cate, and Ben M Harvey. Numerosity tuning in human association cortices and local image contrast representations in early visual cortex. *Nature Communications*, 13(1):1340, 2022.
- [39] Andreas Nieder. The neuronal code for number. *Nature Reviews Neuroscience*, 17(6):366–382, 2016.
- [40] Tommaso Boccato, Alberto Testolin, and Marco Zorzi. Learning numerosity representations with transformers: Number generation tasks and out-of-distribution generalization. *Entropy*, 23(7):857, 2021.
- [41] Alberto Testolin, Will Y Zou, and James L McClelland. Numerosity discrimination in deep neural networks: Initial competence, developmental refinement and experience statistics. *Developmental science*, 23(5):e12940, 2020.
- [42] Eleanor Mundy and Camilla K Gilmore. Children’s mapping between symbolic and nonsymbolic representations of number. *Journal of experimental child psychology*, 103(4):490–502, 2009.
- [43] Iddo Drori, Sarah Zhang, Reece Shuttleworth, Leonard Tang, Albert Lu, Elizabeth Ke, Kevin Liu, Linda Chen, Sunny Tran, Newman Cheng, et al. A neural network solves, explains, and generates university math problems by program synthesis and few-shot learning at human level. *Proceedings of the National Academy of Sciences*, 119(32):e2123433119, 2022.
- [44] Alexis Palmer, Noah A Smith, and Arthur Spirling. Using proprietary language models in academic research requires explicit justification. *Nature Computational Science*, 4(1):2–3, 2024.
- [45] Alberto Testolin. The challenge of modeling the acquisition of mathematical concepts. *Frontiers in human neuroscience*, 14:100, 2020.
- [46] Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.
- [47] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022.