

ViT (Vision Transformers)

[ViT paper](#)

[Video explanation](#)

Best Transformer intro: [https://jalammar.github.io/illustrate...](https://jalammar.github.io/illustrating-vision-transformer/)

CNNs vs ViT: <https://arxiv.org/abs/2108.08810>

CNNs vs ViT Blog: [https://towardsdatascience.com/do-vis...](https://towardsdatascience.com/do-vision-transformers-really-work-4f3a2a2a2a2a)

Swin Transformer: <https://arxiv.org/abs/2103.14030>

DeiT: <https://arxiv.org/abs/2012.12877>

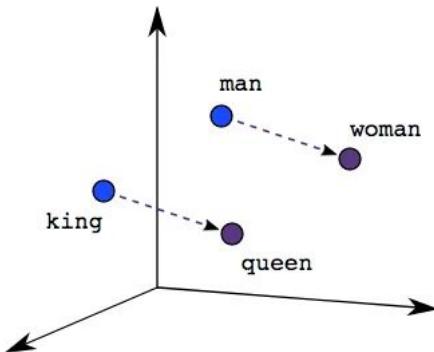
Visualizations:

<https://jacobgil.github.io/deeplearning/vision-transformer-explainability>

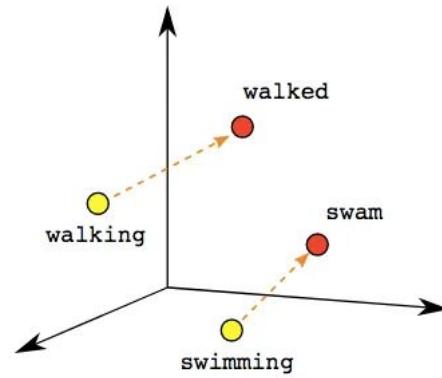
<https://arxiv.org/pdf/2012.09838v1.pdf>

Representación semántica en espacio vectorial

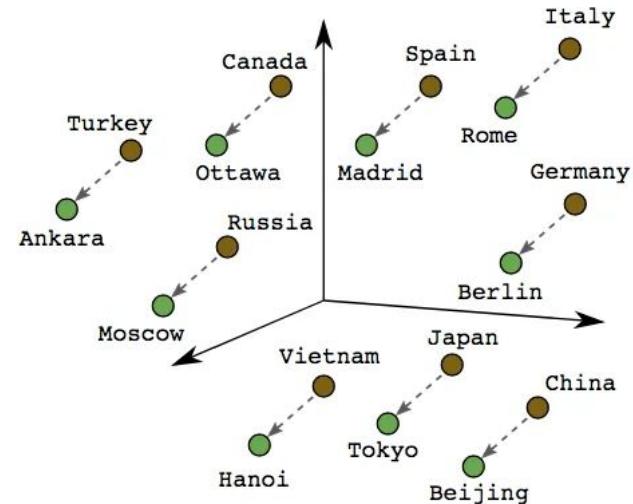
Token embeddings



Male-Female

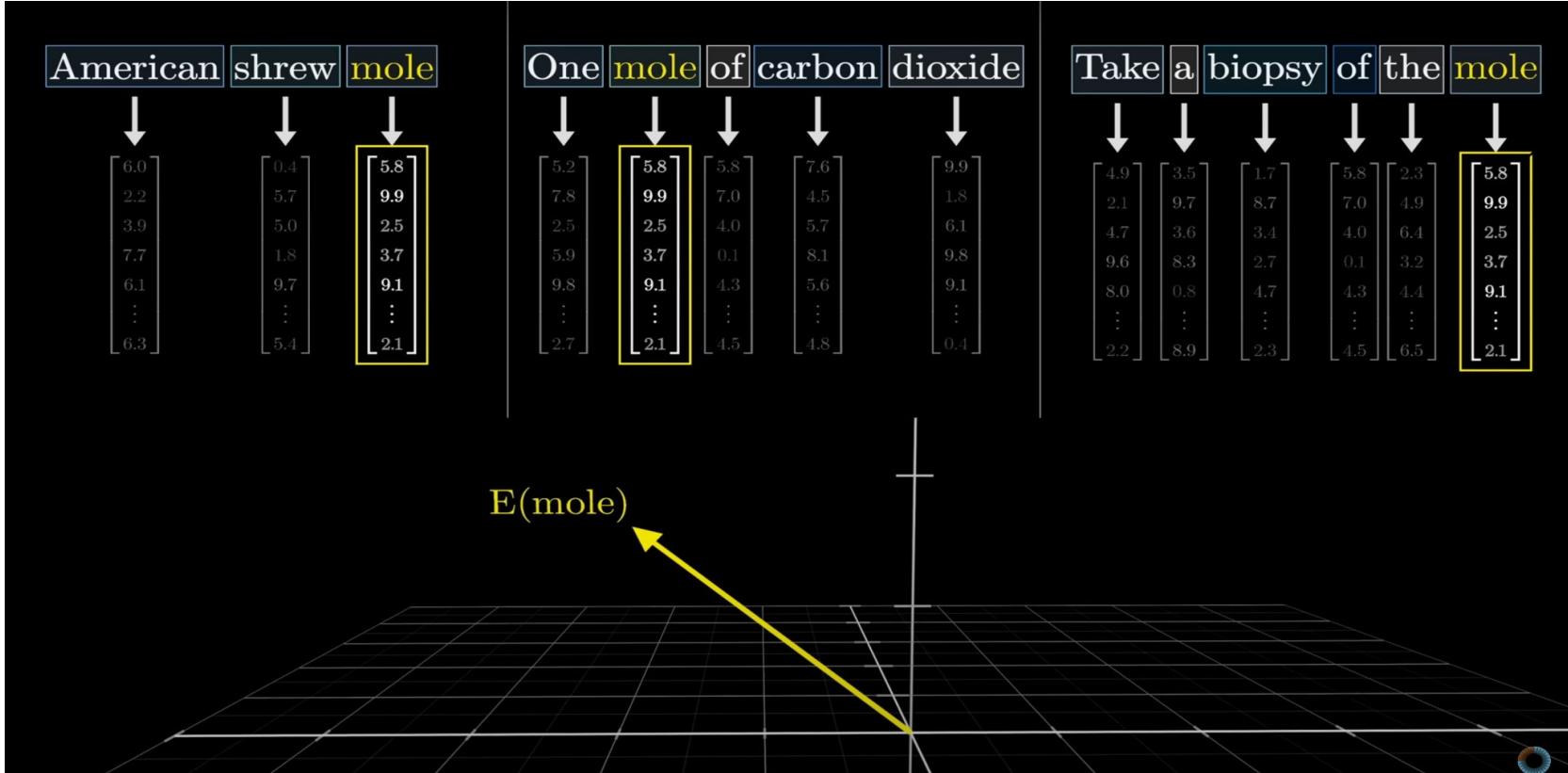


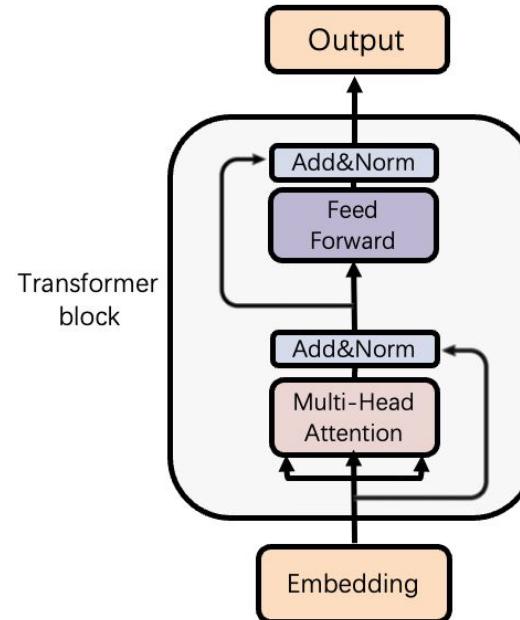
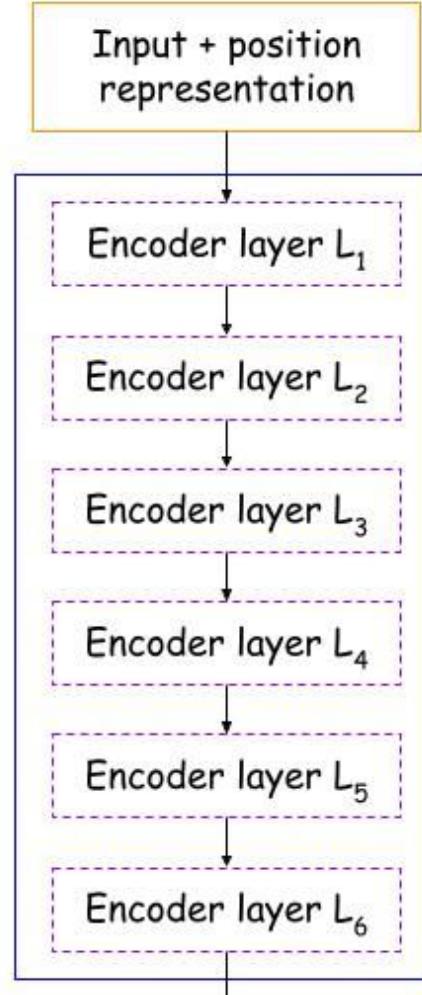
Verb Tense



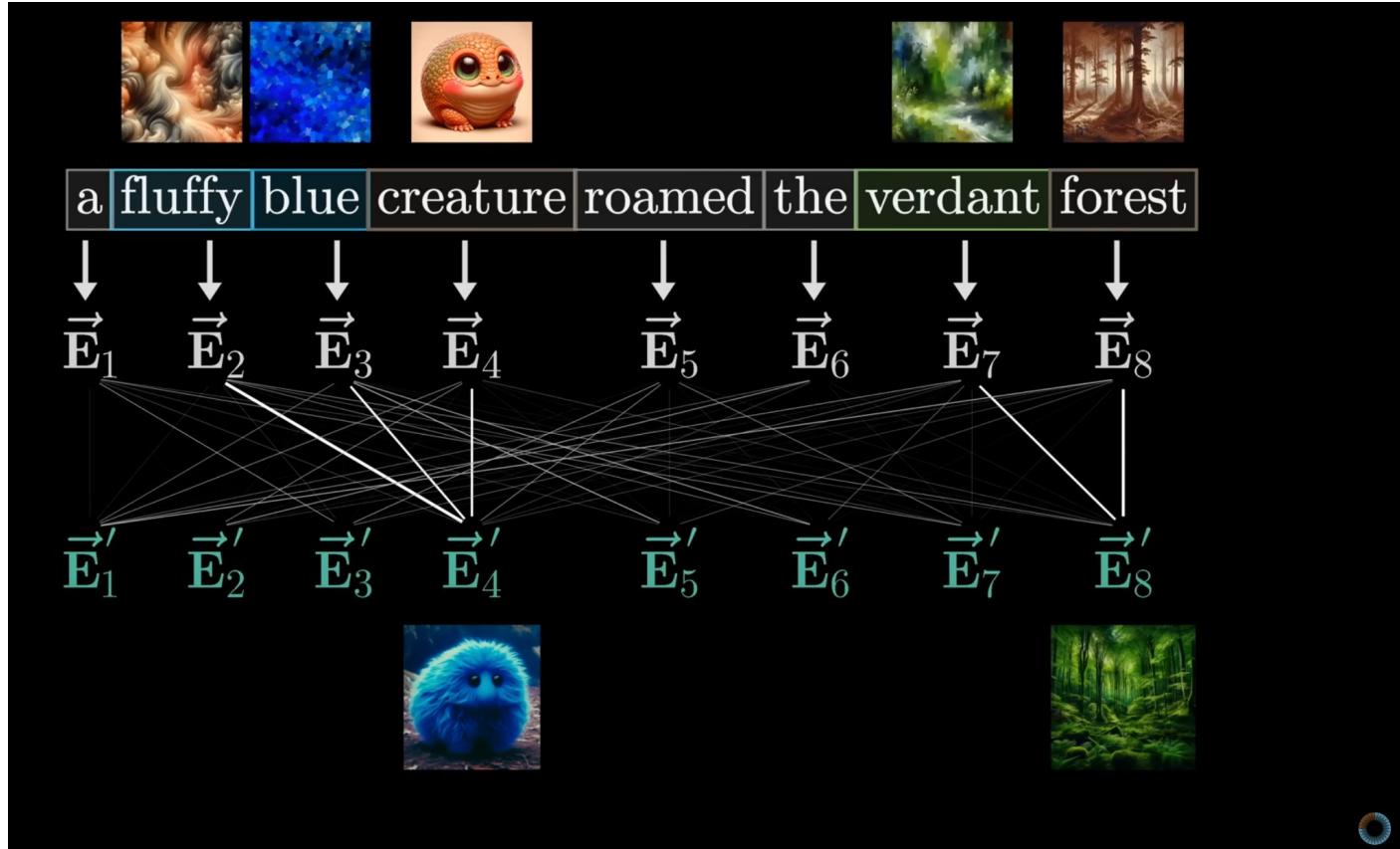
Country-Capital

Non contextual embeddings ([3blue1brown](#))



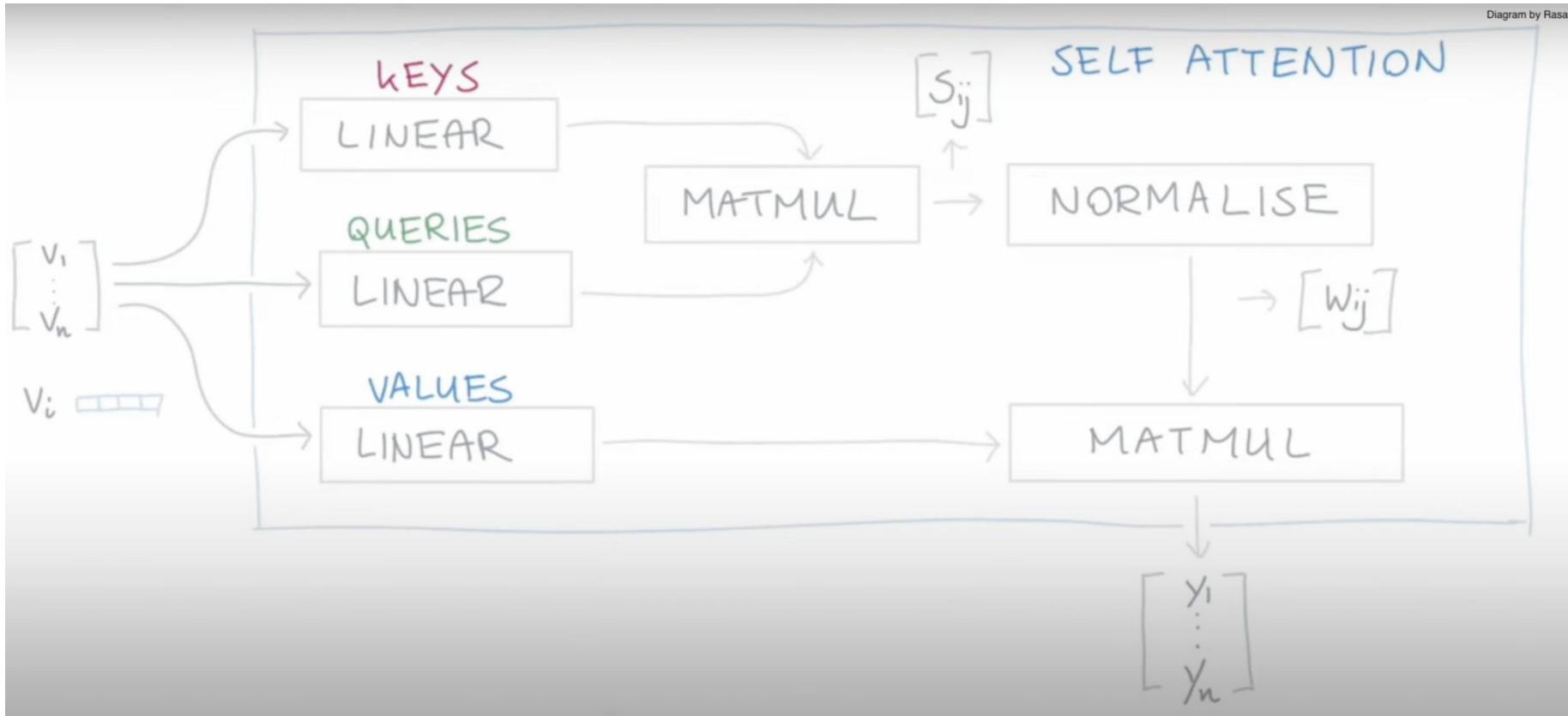


Contextual embeddings (Transformer architecture)



$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^\top}{\sqrt{d_k}} \right) V$$

Self Attention



Concept	CNN	Self-Attention
Input	Image: $H \times W \times C$	Sequence: $T \times d_{\text{model}}$
"Filters"	Learnable kernels (e.g., $3 \times 3 \times C$)	Learnable linear projections: W^Q, W^K, W^V
Computation	Dot product of kernel with input patch	Dot product of query with key: QK^\top
Output	Feature maps per filter	Weighted sum of values: $\text{softmax}(QK^\top)V$
Spatial Context	Local (e.g., 3×3 window)	Global (attends to all tokens)
Channels	Depth (number of filters)	Embedding dimension (model depth)

CNNs apply fixed-size **local filters** to detect patterns in space.

Attention uses **adaptive, content-dependent filters** (via dot products) that can look at the **whole sequence**

Analogía canales / embedding / features

Contexto 2021

Transformers aplicados a NLP

- pre entrenamiento con corpus grandes
- ajuste fino con dataset en tareas específicas
- modelos y datasets en crecimiento, sin indicios de saturación de performance

Computer Vision:

- Estado del arte para CNNs
- Se empezaron a combinar con self attention

Resultados preliminares con transformers

Dataset medianos como ImageNet:

- Pocos puntos debajo de ResNets de tamaños comparables
- Resultado esperado
- Carecen del sesgo inductivo de las CNNs (equivarianza a la translación y localidad)
- Necesitan mucha data para generalizar (debido a no tener sesgos)

Dataset grandes (14M a 300M)

- ImageNet-21k: 21k classes and 14M images
- JFT: 18k classes and 303M high-resolution images

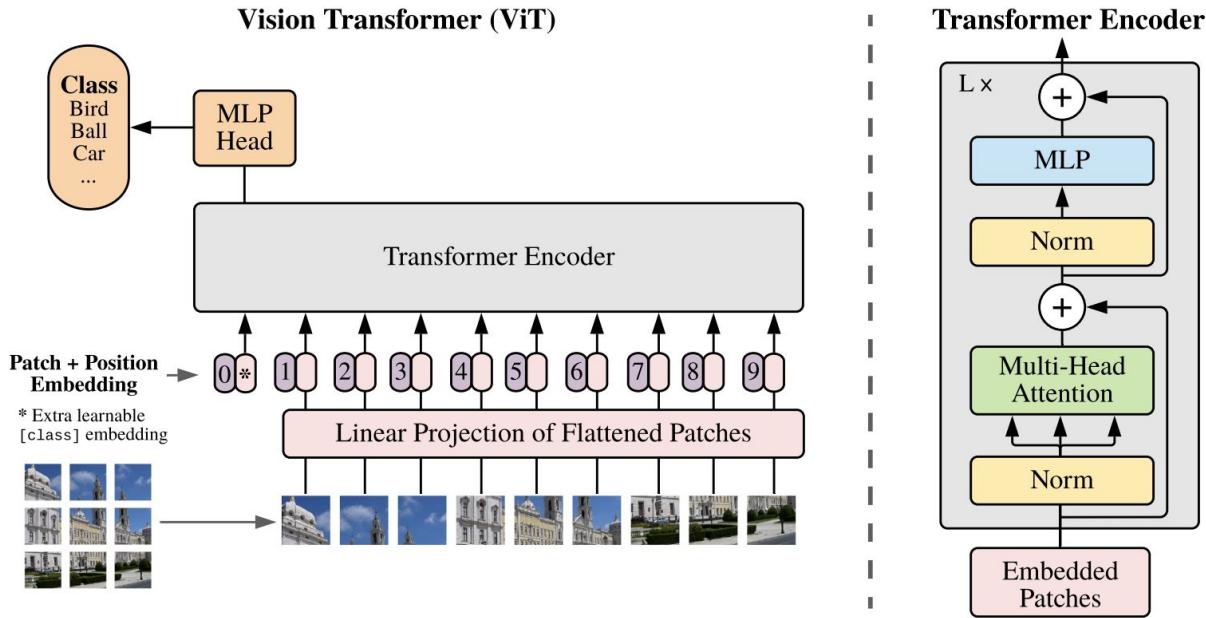


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

Model	Layers	Hidden size D	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

Table 1: Details of Vision Transformer model variants.

	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21k (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	88.55 ± 0.04	87.76 ± 0.03	85.30 ± 0.02	87.54 ± 0.02	88.4/88.5*
ImageNet ReaL	90.72 ± 0.05	90.54 ± 0.03	88.62 ± 0.05	90.54	90.55
CIFAR-10	99.50 ± 0.06	99.42 ± 0.03	99.15 ± 0.03	99.37 ± 0.06	—
CIFAR-100	94.55 ± 0.04	93.90 ± 0.05	93.25 ± 0.05	93.51 ± 0.08	—
Oxford-IIIT Pets	97.56 ± 0.03	97.32 ± 0.11	94.67 ± 0.15	96.62 ± 0.23	—
Oxford Flowers-102	99.68 ± 0.02	99.74 ± 0.00	99.61 ± 0.02	99.63 ± 0.03	—
VTAB (19 tasks)	77.63 ± 0.23	76.28 ± 0.46	72.72 ± 0.21	76.29 ± 1.70	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k

Table 2: Comparison with state of the art on popular image classification benchmarks. We report mean and standard deviation of the accuracies, averaged over three fine-tuning runs. Vision Transformer models pre-trained on the JFT-300M dataset outperform ResNet-based baselines on all datasets, while taking substantially less computational resources to pre-train. ViT pre-trained on the smaller public ImageNet-21k dataset performs well too. *Slightly improved 88.5% result reported in Touvron et al. (2020).

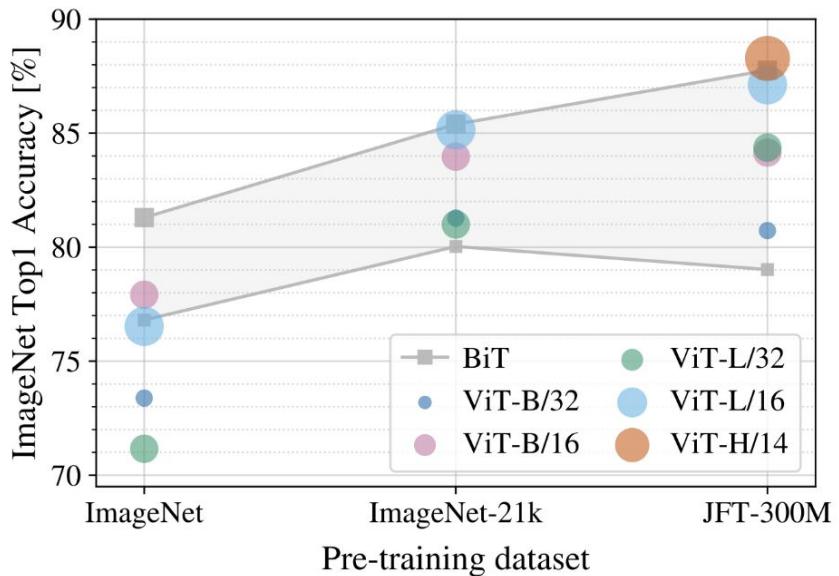


Figure 3: Transfer to ImageNet. While large ViT models perform worse than BiT ResNets (shaded area) when pre-trained on small datasets, they shine when pre-trained on larger datasets. Similarly, larger ViT variants overtake smaller ones as the dataset grows.

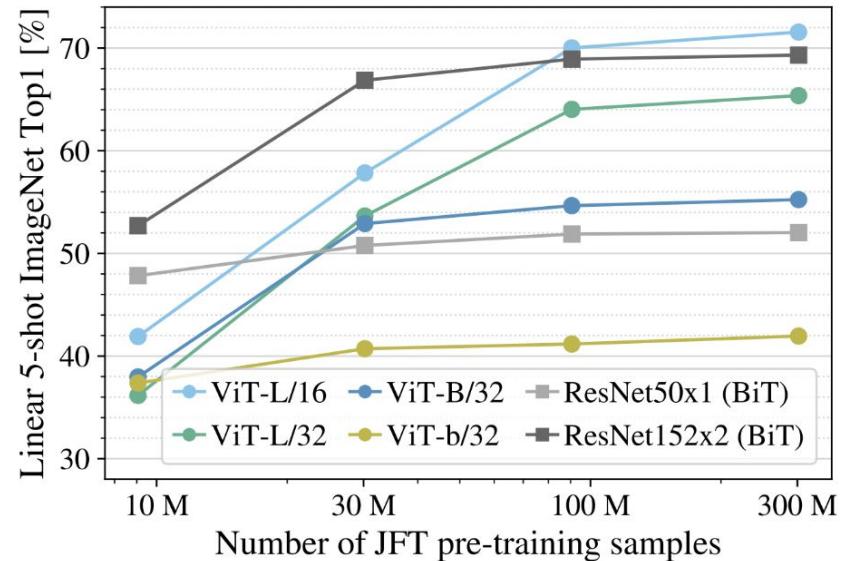


Figure 4: Linear few-shot evaluation on ImageNet versus pre-training size. ResNets perform better with smaller pre-training datasets but plateau sooner than ViT, which performs better with larger pre-training. ViT-b is ViT-B with all hidden dimensions halved.

Training

Models	Dataset	Epochs	Base LR	LR decay	Weight decay	Dropout
ViT-B/{16,32}	JFT-300M	7	$8 \cdot 10^{-4}$	linear	0.1	0.0
ViT-L/32	JFT-300M	7	$6 \cdot 10^{-4}$	linear	0.1	0.0
ViT-L/16	JFT-300M	7/14	$4 \cdot 10^{-4}$	linear	0.1	0.0
ViT-H/14	JFT-300M	14	$3 \cdot 10^{-4}$	linear	0.1	0.0
R50x{1,2}	JFT-300M	7	10^{-3}	linear	0.1	0.0
R101x1	JFT-300M	7	$8 \cdot 10^{-4}$	linear	0.1	0.0
R152x{1,2}	JFT-300M	7	$6 \cdot 10^{-4}$	linear	0.1	0.0
R50+ViT-B/{16,32}	JFT-300M	7	$8 \cdot 10^{-4}$	linear	0.1	0.0
R50+ViT-L/32	JFT-300M	7	$2 \cdot 10^{-4}$	linear	0.1	0.0
R50+ViT-L/16	JFT-300M	7/14	$4 \cdot 10^{-4}$	linear	0.1	0.0
ViT-B/{16,32}	ImageNet-21k	90	10^{-3}	linear	0.03	0.1
ViT-L/{16,32}	ImageNet-21k	30/90	10^{-3}	linear	0.03	0.1
ViT-*	ImageNet	300	$3 \cdot 10^{-3}$	cosine	0.3	0.1

Table 3: Hyperparameters for training. All models are trained with a batch size of 4096 and learning rate warmup of 10k steps. For ImageNet we found it beneficial to additionally apply gradient clipping at global norm 1. Training resolution is 224.

Fine tuning

- SGD con momentum = 0.5
- Small Grid Search para learning rates
- resolución: 384 (diferente que train - práctica común)
- batch_size = 512 (diferente que train)
- no weight decay
- Se reemplaza la Head (dos capas densas) por una de Dx C donde D es la dimensión de los embeddings y C la cantidad de clases (inicializada en 0)

Algunas cosas cambian según el dataset

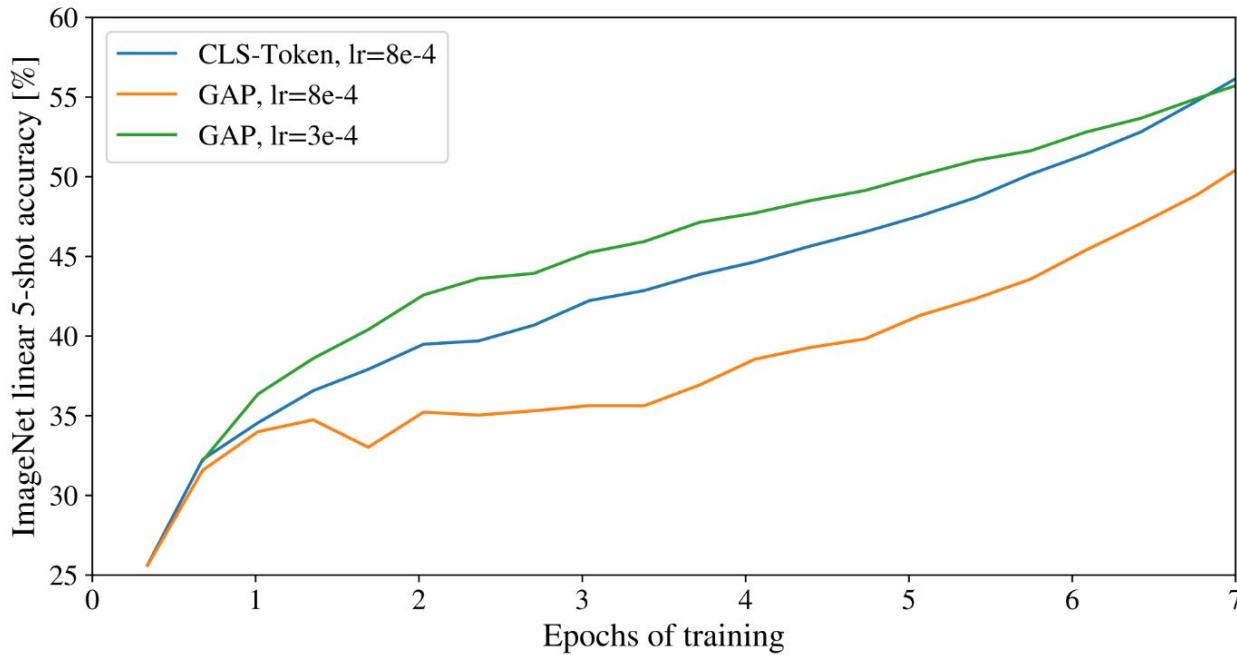


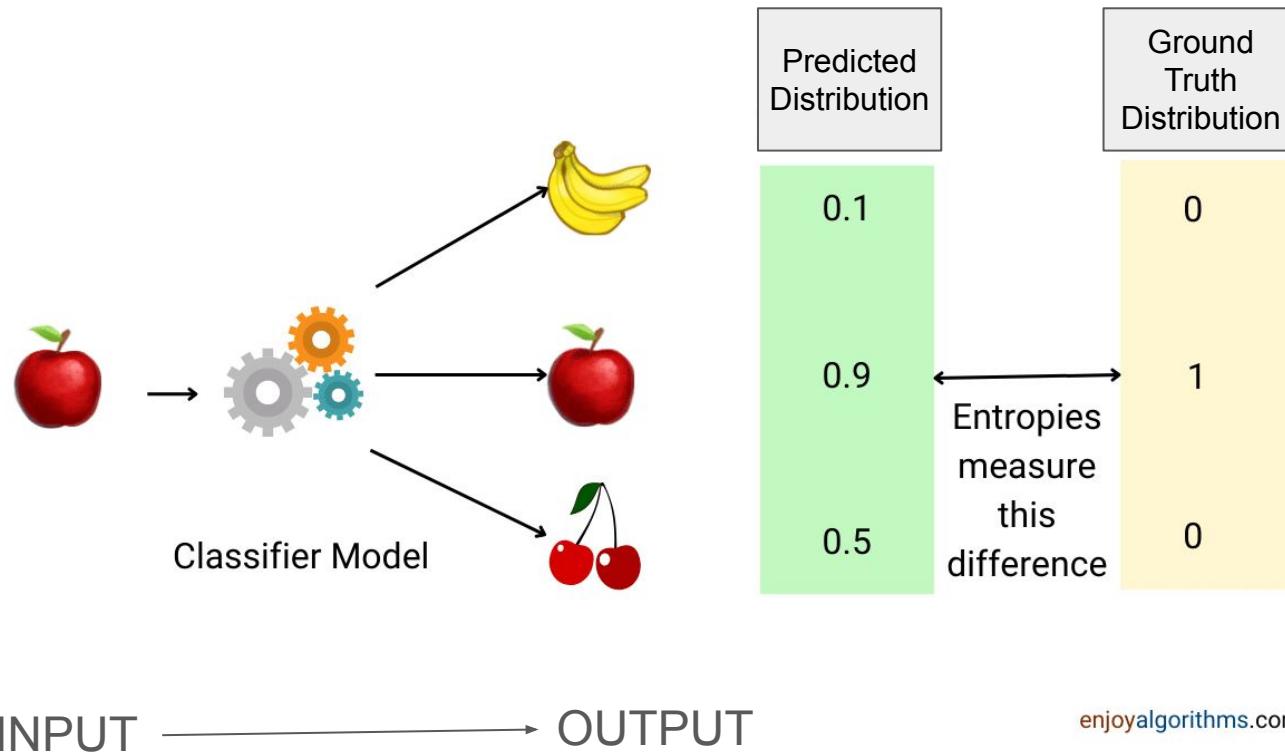
Figure 9: Comparison of class-token and global average pooling classifiers. Both work similarly well, but require different learning-rates.

Foundational Models

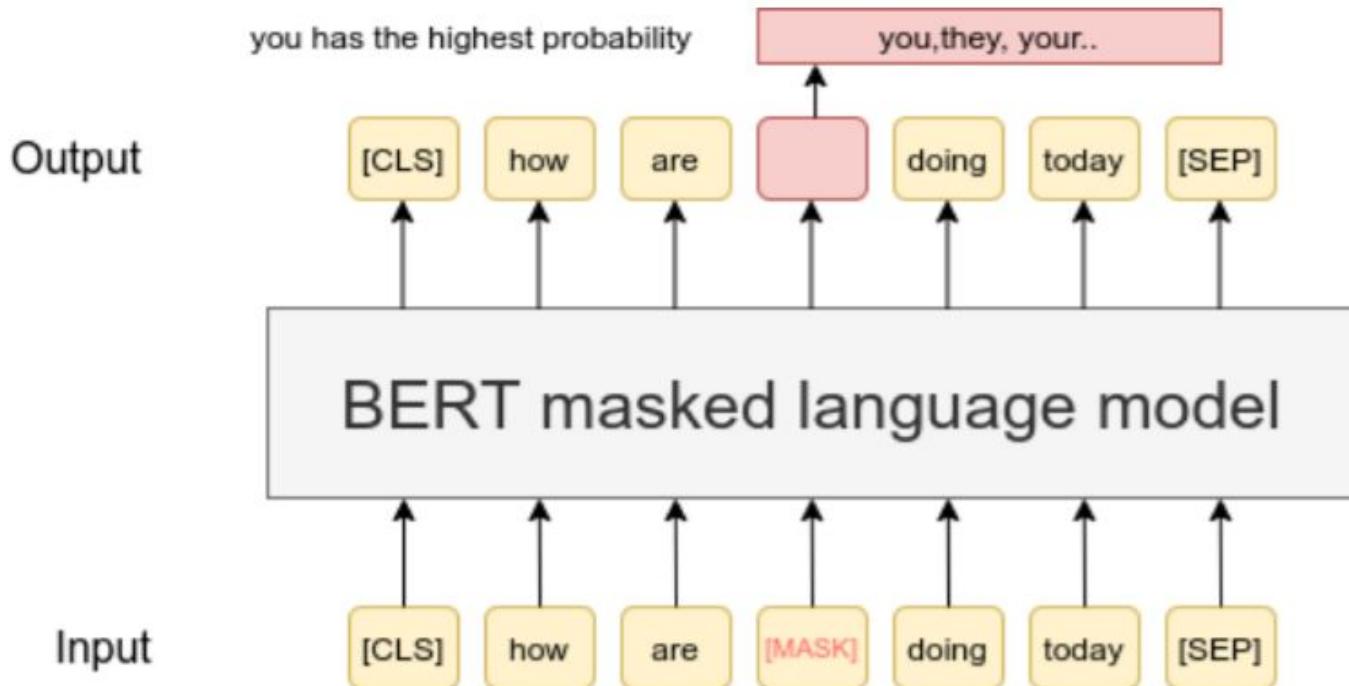
AI Foundational model

1. **Pre-trained on Large Datasets** 
 - a. Self-supervised (unsupervised)
 - b. Massive dataset
2. Generalizable Representations
 - a. Learn Meaningful features (semantics)
 - b. Embeddings
3. Transferable
 - a. Fine-tunable in small diverse datasets
 - b. Multiple downstream tasks: Classification, Generation, Prediction

Supervised learning Training - Needs labeling



Self-supervised (Encoder) - BERT ([Paper](#))



you	0.459
we	0.212
they	0.153
things	0.038
people	0.035
others	0.019
students	0.018
kids	0.016

Bidirectional Encoder Representations Transformer

Randomly masks 15% of the words.

Self-supervised (decoder) - GPT

Inputs

Trans | formers | are | great | for | machine | learning | !



GPT (Generative Pre-Train Transformer) Model

Trans | formers | are | great | for | machine | learning | ! Labels

Inference - GPT

My



name

My

name



is

My

name

is



Sylvain

My

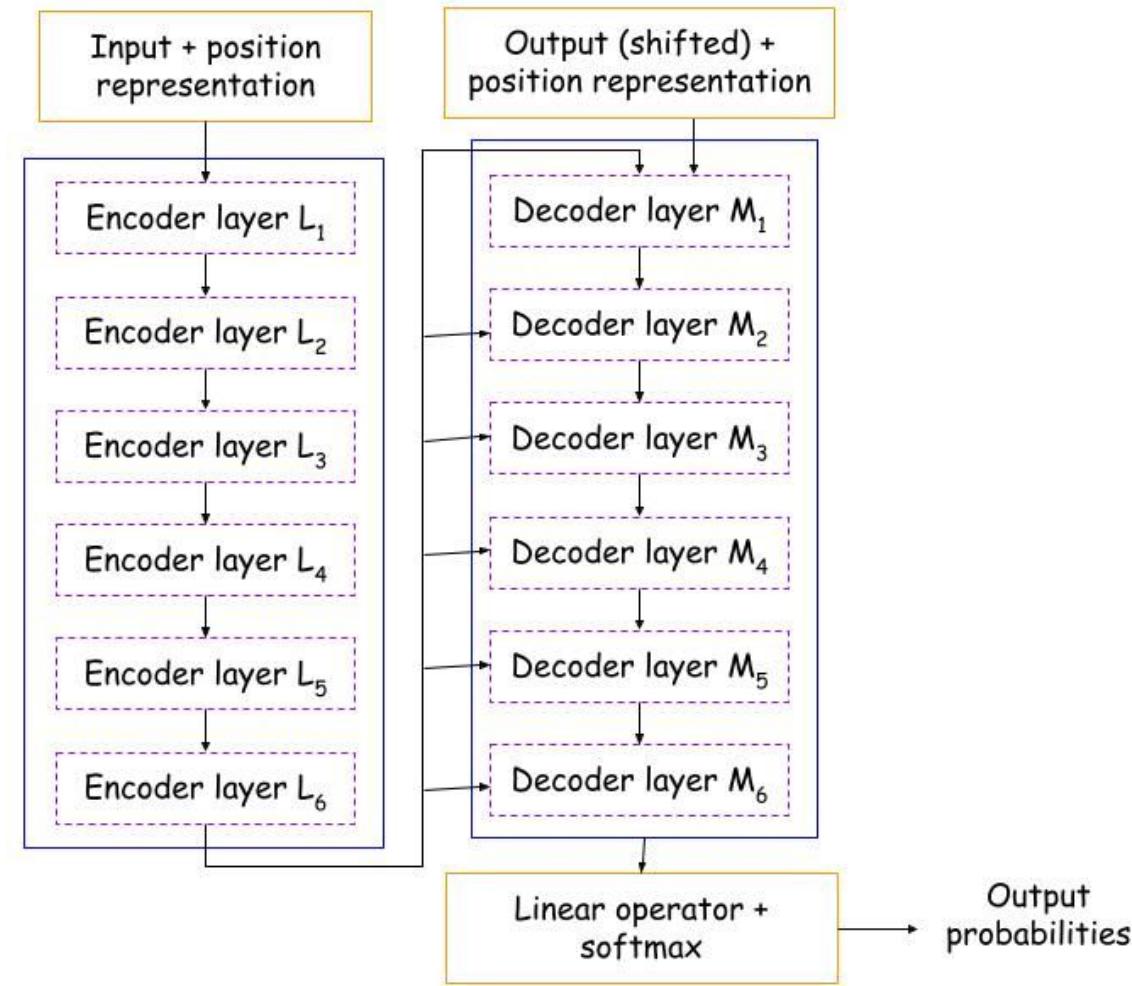
name

is

Sylvain



.



What does the model learn?

Sintaxis: Language structure, grammar rules, sentence construction, etc. Learned not by rules but by exposure to billions of examples

Semantics: word meanings, concept relationships, and how words/phrases are used in different contexts.

Context: They pick up on contextual dependencies, even across long passages.

Zero-shot

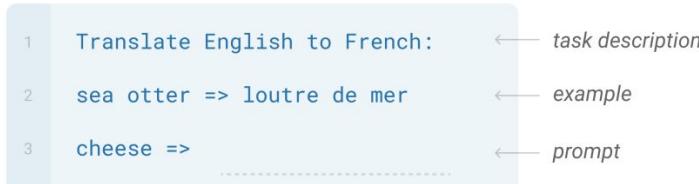
Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



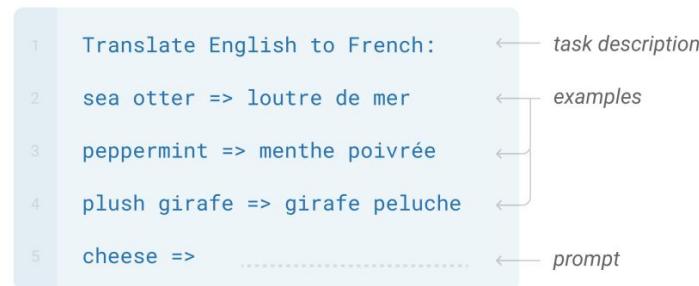
One-shot

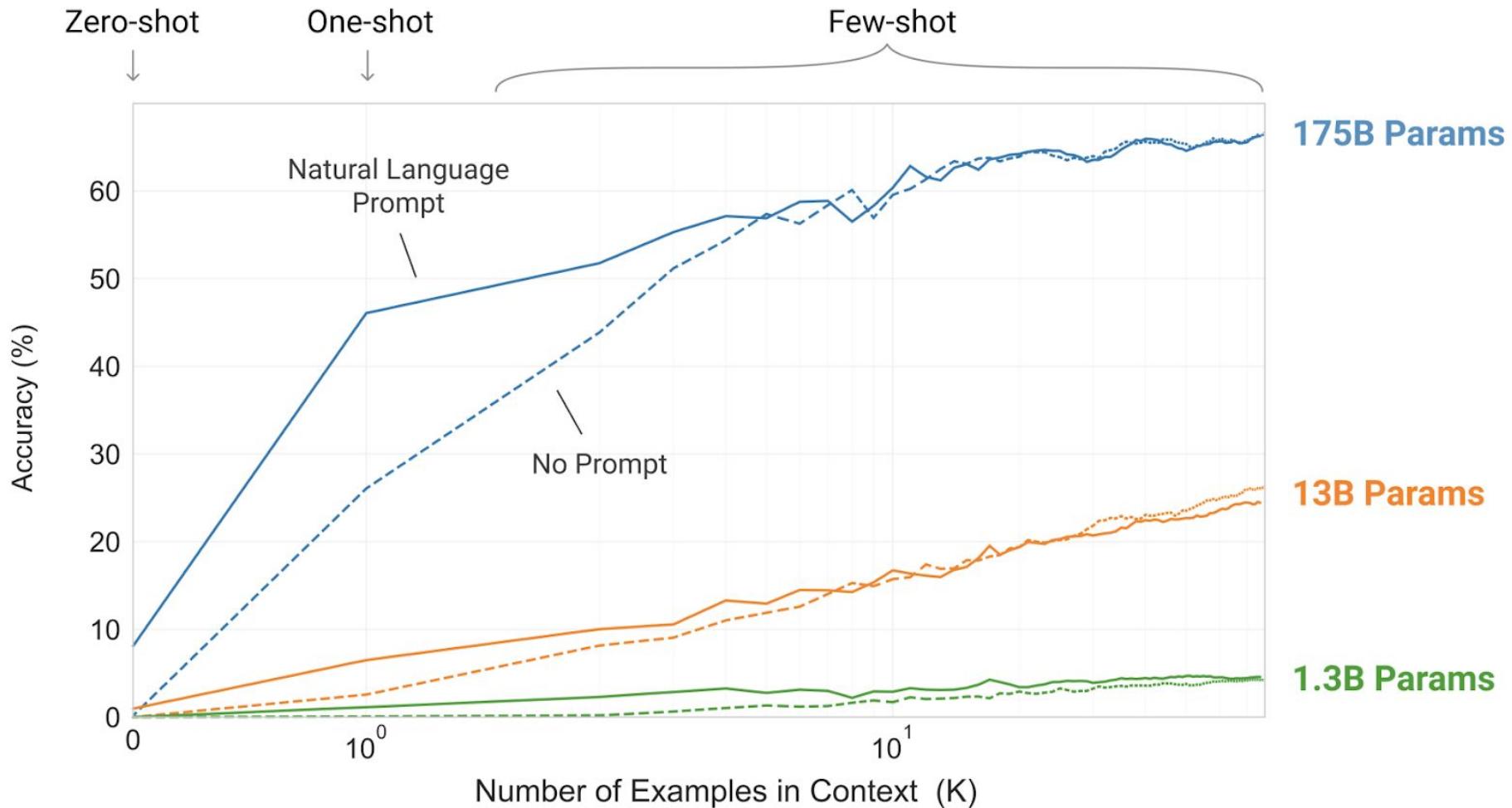
In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.





Massive dataset

Token: Word fraction, example: tele-vision

Total tokens: 50k a 200K

(GPT3: 50,257, LLaMA3: 128,256)

Training Total tokens: 300B a 15T

(LLaMA1: 1.4T, LLaMA3: 15T)

Total different data tokens:

(LLaMA1: <1T, LLaMA3: ~3/5T)

More than one epoch per training

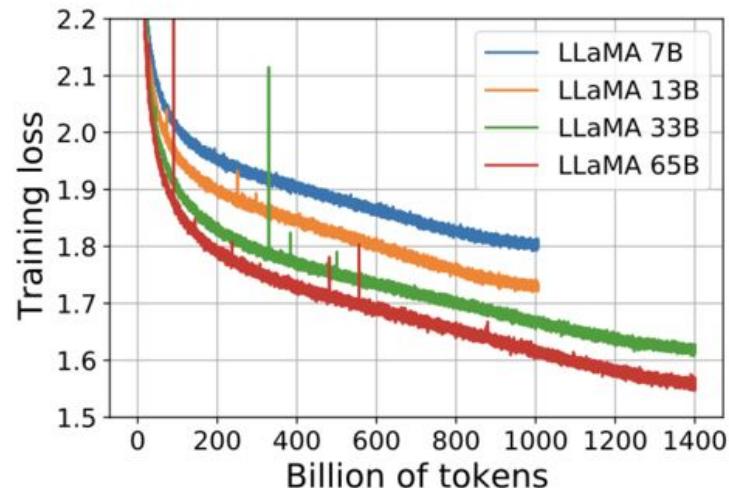


Figure 1: **Training loss over train tokens for the 7B, 13B, 33B, and 65 models.** LLaMA-33B and LLaMA-65B were trained on 1.4T tokens. The smaller models were trained on 1.0T tokens. All models are trained with a batch size of 4M tokens.

AI Foundational model

1. Pre-trained on Large Datasets
 - a. Self-supervised (unsupervised)
 - b. Massive dataset
2. Generalizable Representations
 - a. Learn Meaningful features (semantics)
 - b. Embeddings
3. Transferable
 - a. Fine-tunable in small diverse datasets
 - b. Multiple downstream tasks: Classification, Generation, Prediction

AI Foundational model

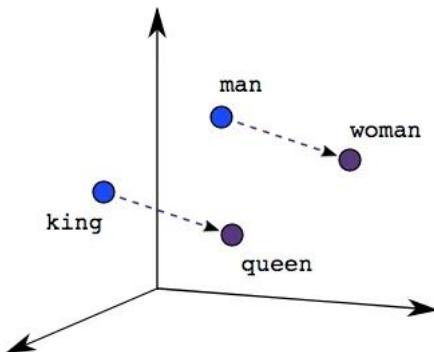
1. Pre-trained on Large Datasets
 - a. Self-supervised (unsupervised)
 - b. Massive dataset
2. **Generalizable Representations** 

 - a. Learn Meaningful features (semantics)
 - b. Embeddings

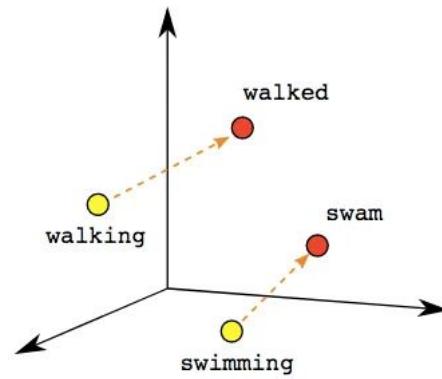
3. Transferable
 - a. Fine-tunable in small diverse datasets
 - b. Multiple downstream tasks: Classification, Generation, Prediction

Representación semántica en espacio vectorial

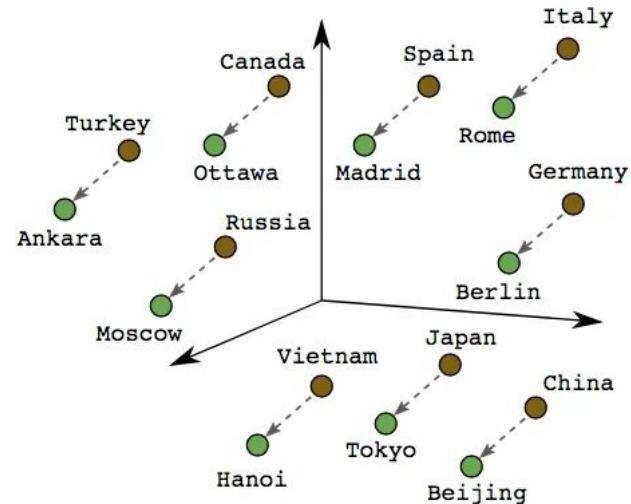
Token embeddings



Male-Female

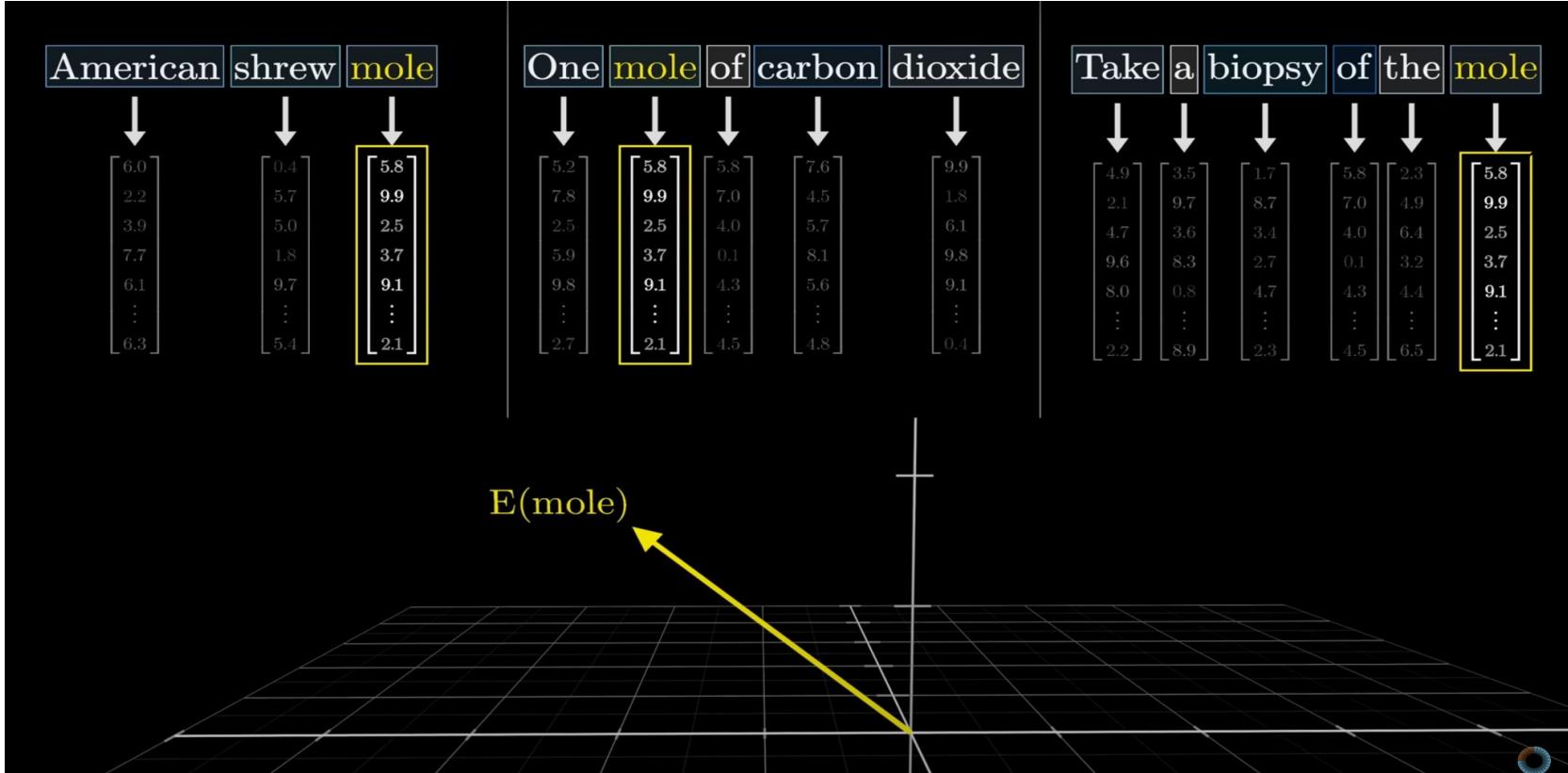


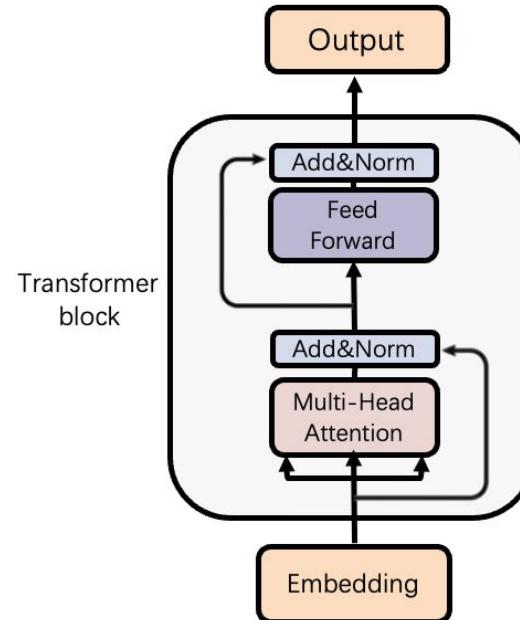
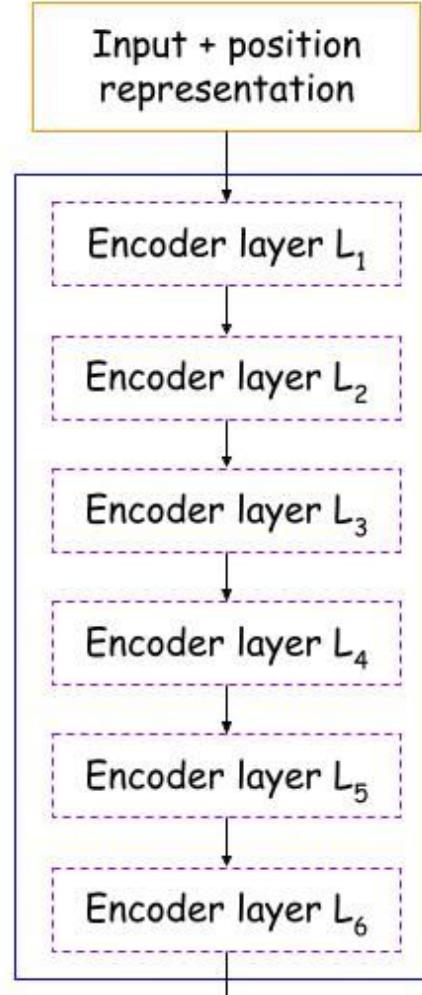
Verb Tense



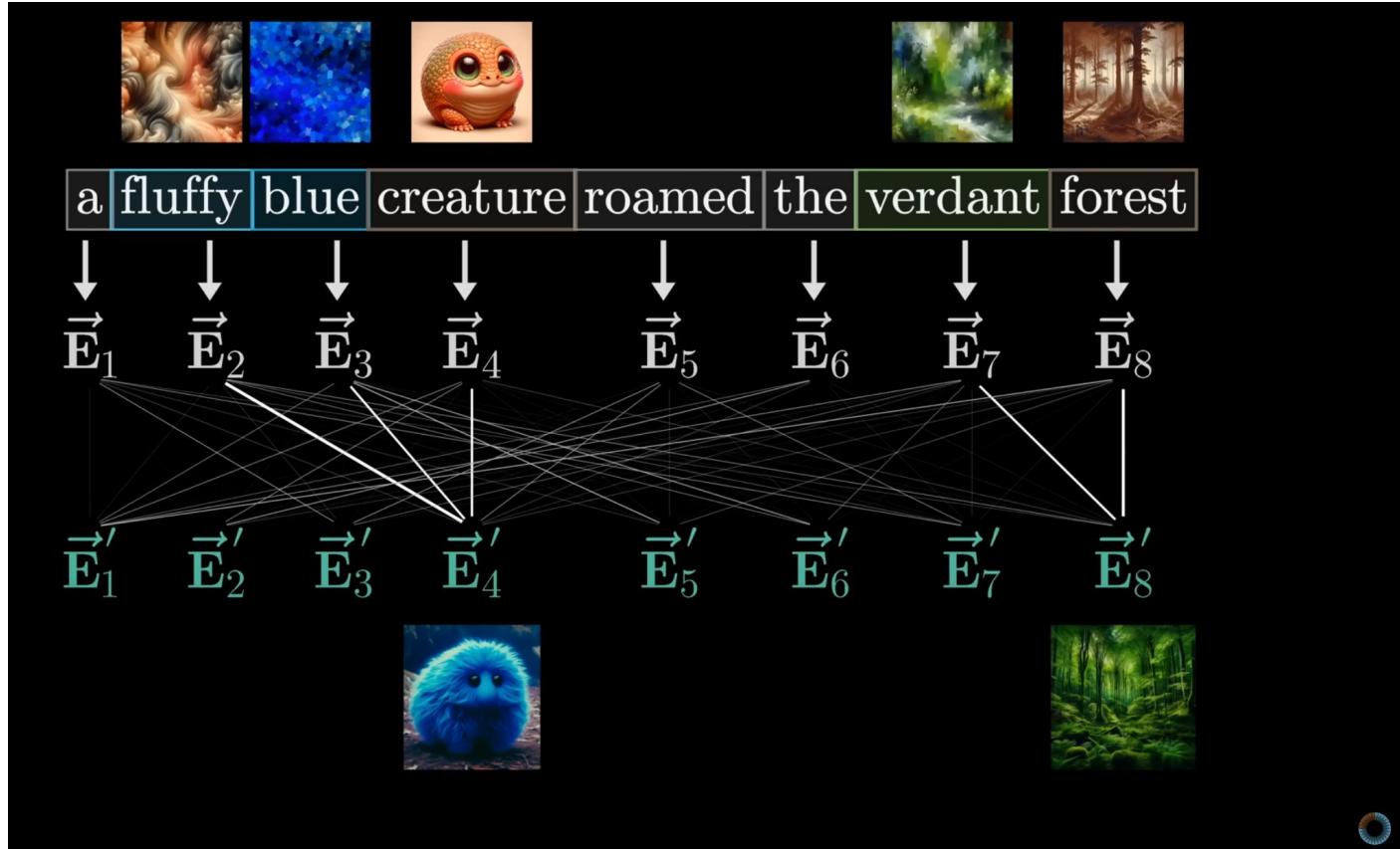
Country-Capital

Non contextual embeddings ([3blue1brown](#))

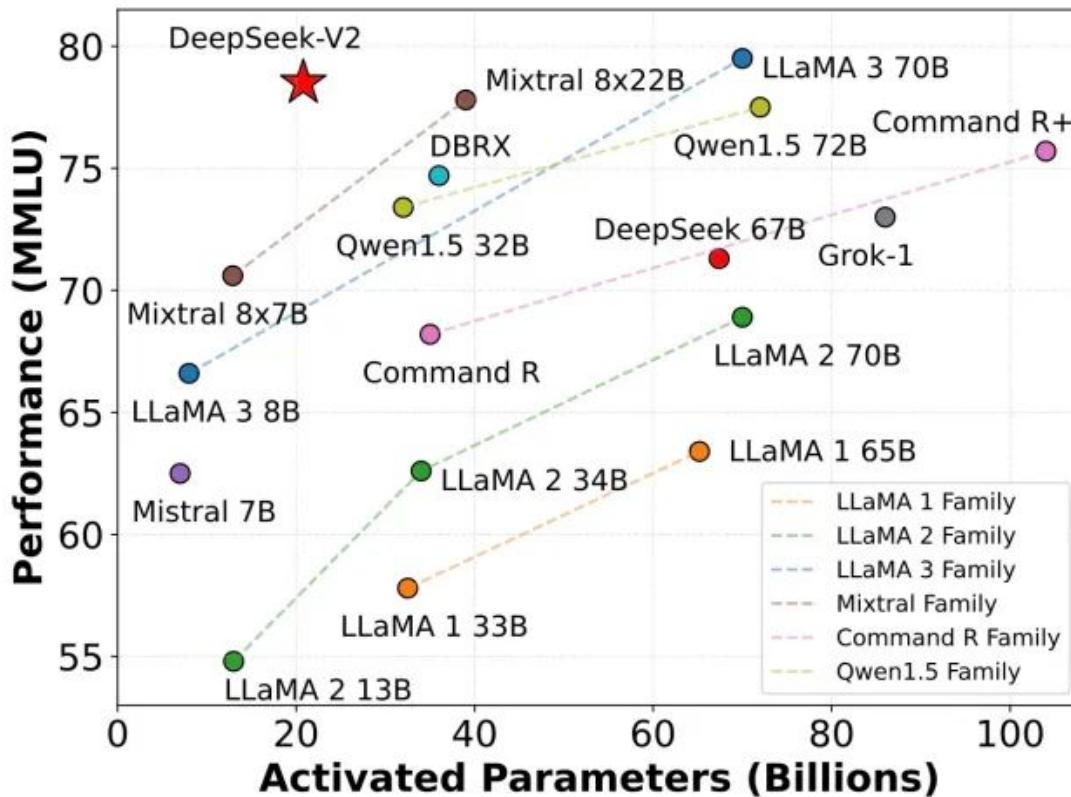




Contextual embeddings (Transformer architecture)



Will this tendency continue?



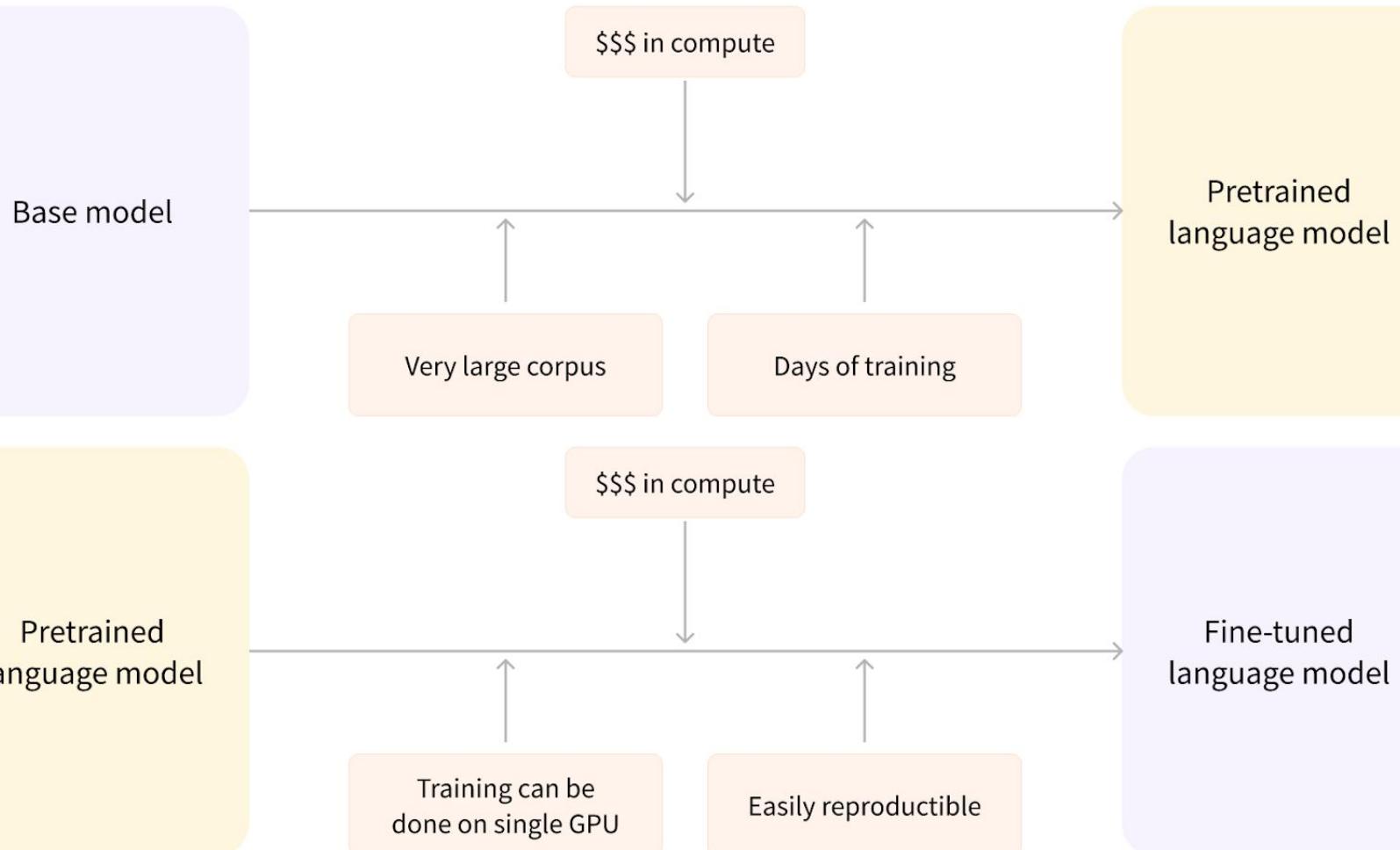
Parameters	Tokens
400 Million	8.0 Billion
1 Billion	20.2 Billion
10 Billion	205.1 Billion
67 Billion	1.5 Trillion
175 Billion	3.7 Trillion
280 Billion	5.9 Trillion
520 Billion	11.0 Trillion
1 Trillion	21.2 Trillion
10 Trillion	216.2 Trillion

AI Foundational model

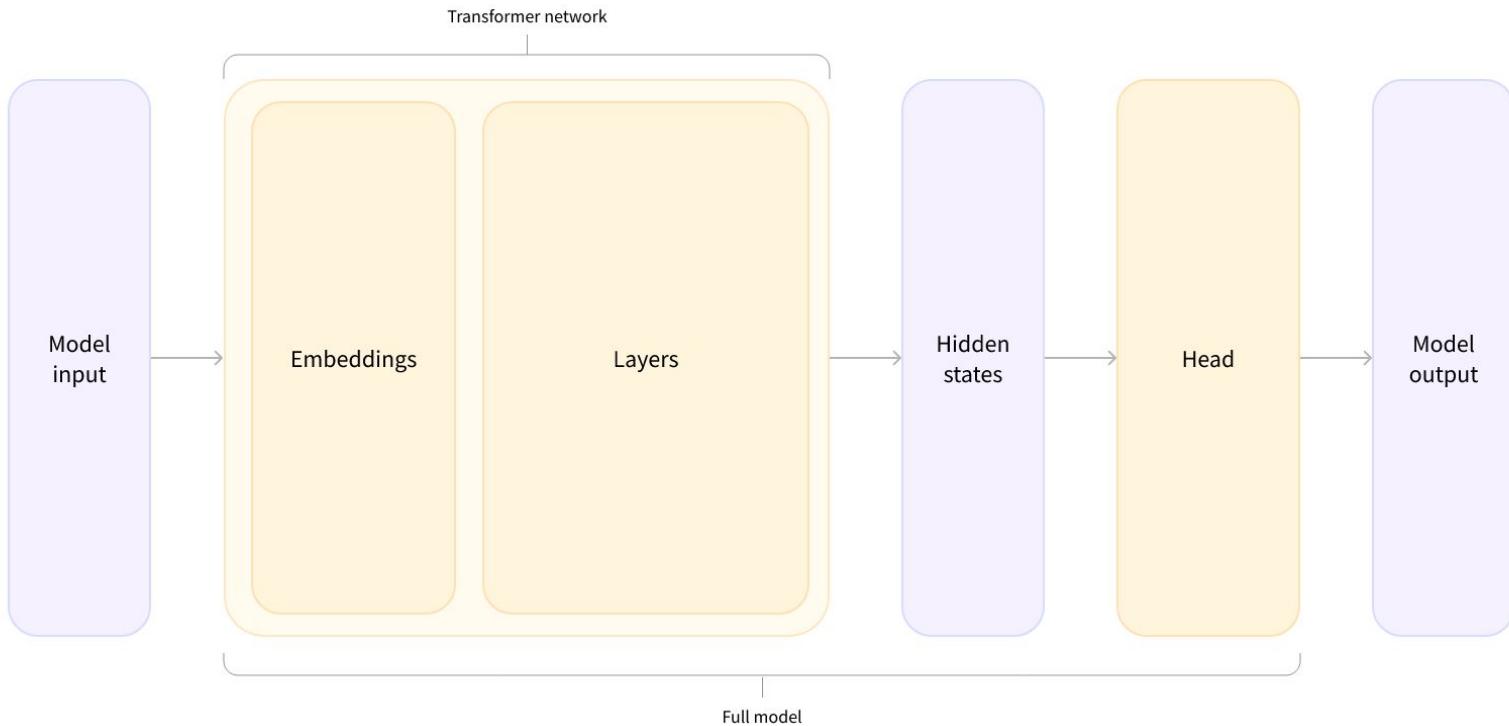
1. Pre-trained on Large Datasets
 - a. Self-supervised (unsupervised)
 - b. Massive dataset
2. Generalizable Representations
 - a. Learn Meaningful features (semantics)
 - b. Embeddings
3. Transferable
 - a. Fine-tunable in small diverse datasets
 - b. Multiple downstream tasks: Classification, Generation, Prediction

AI Foundational model

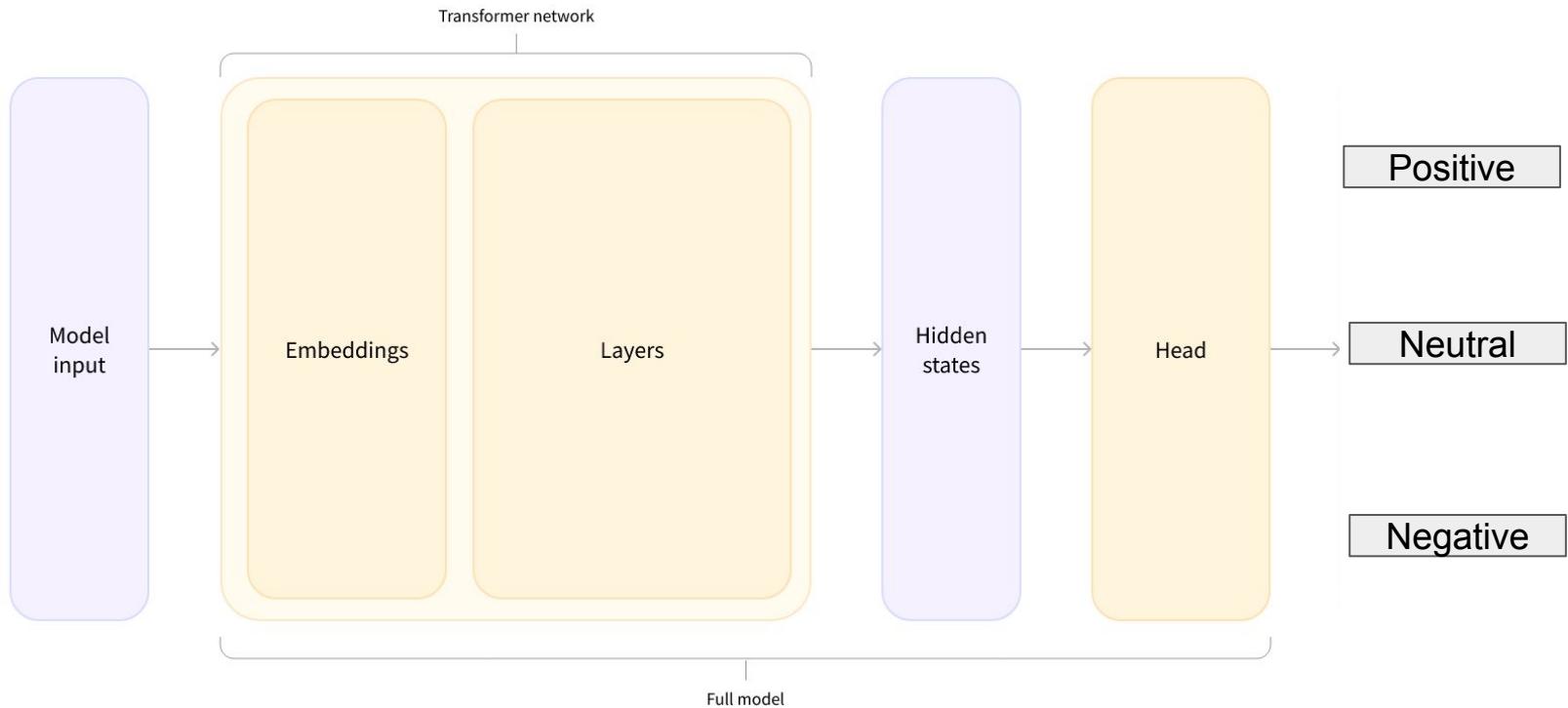
1. Pre-trained on Large Datasets
 - a. Self-supervised (unsupervised)
 - b. Massive dataset
2. Generalizable Representations
 - a. Learn Meaningful features (semantics)
 - b. Embeddings
3. Transferable 
 - a. Fine-tunable in small diverse datasets
 - b. Multiple downstream tasks: Classification, Generation, Regression



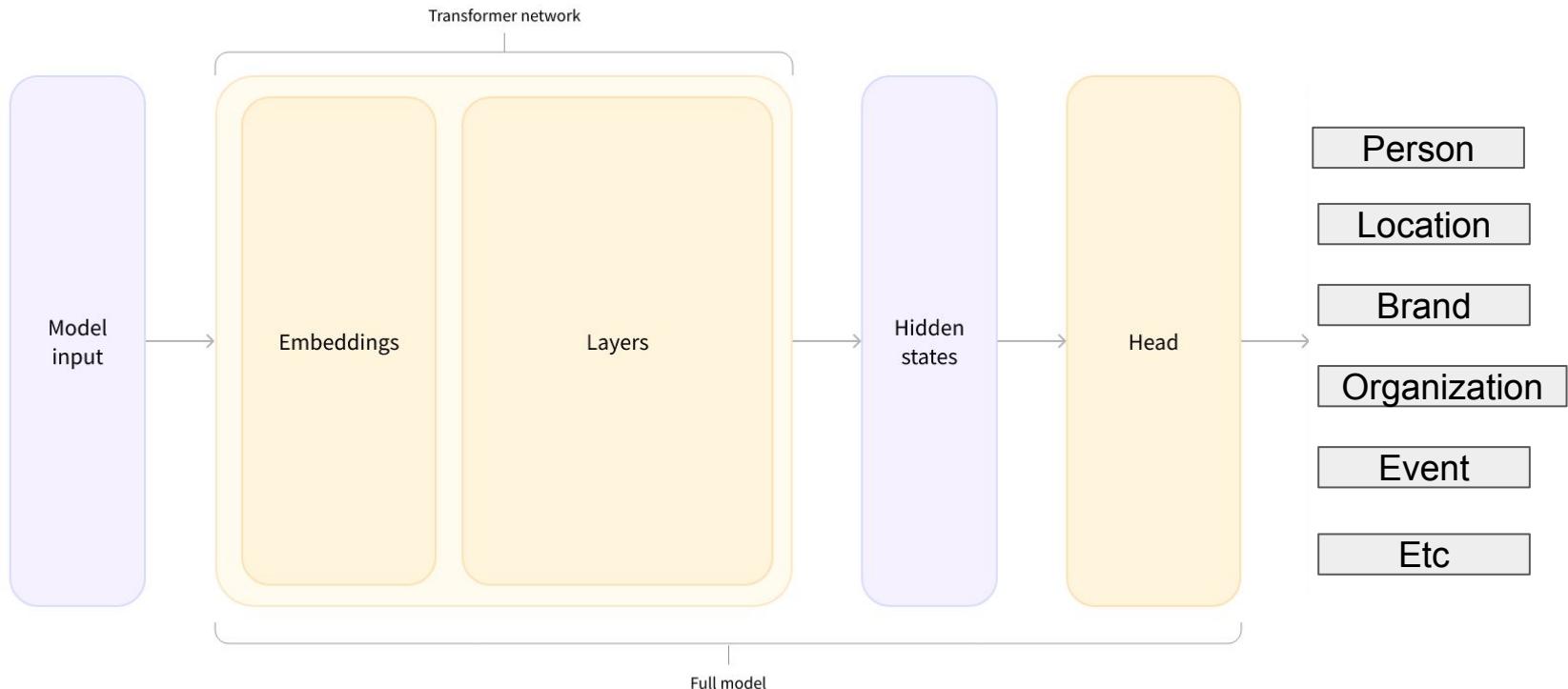
Transferable and fine tuned



Sentiment Analysis (Classification of full sentence)



NER (Named Entity Recognition) | For each token



Generation

 Chat - Instruct

 Translation

 Summarization

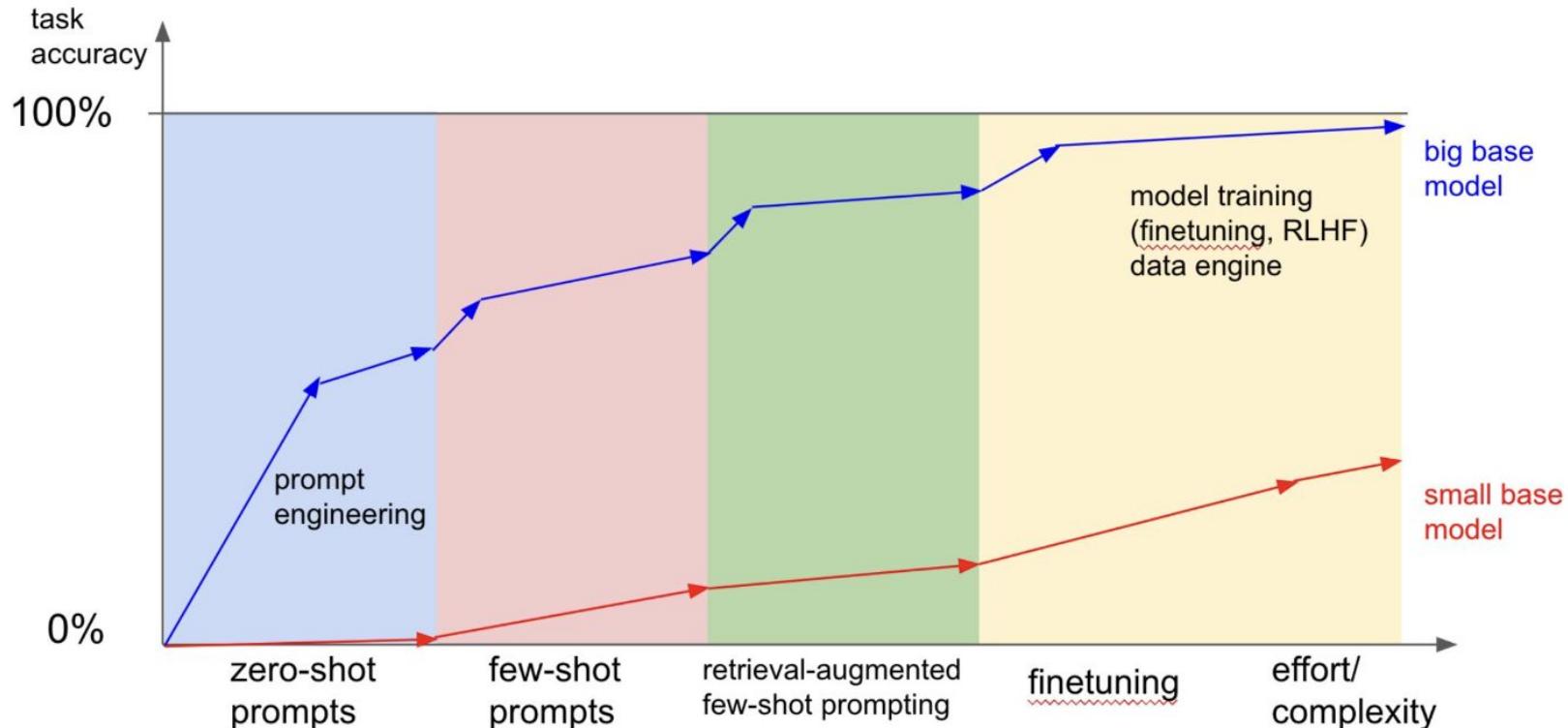
Regression Examples

Stock pricing movement prediction: Use news articles, earnings call transcripts, or SEC filings as input to predict stock returns or volatility.

Credit scoring: Use textual reports (e.g., customer descriptions, loan applications) to regress a risk score

Personality trait estimation: Based on user-written text, LLMs can output Big Five personality trait scores (continuous values).

Mental health screening: Text input (journals, social media) used to predict depression or anxiety scores.



Masked Autoencoders are Scalable Visual Learners

“This paper shows that masked autoencoders (MAE) are scalable self-supervised learners for computer vision”

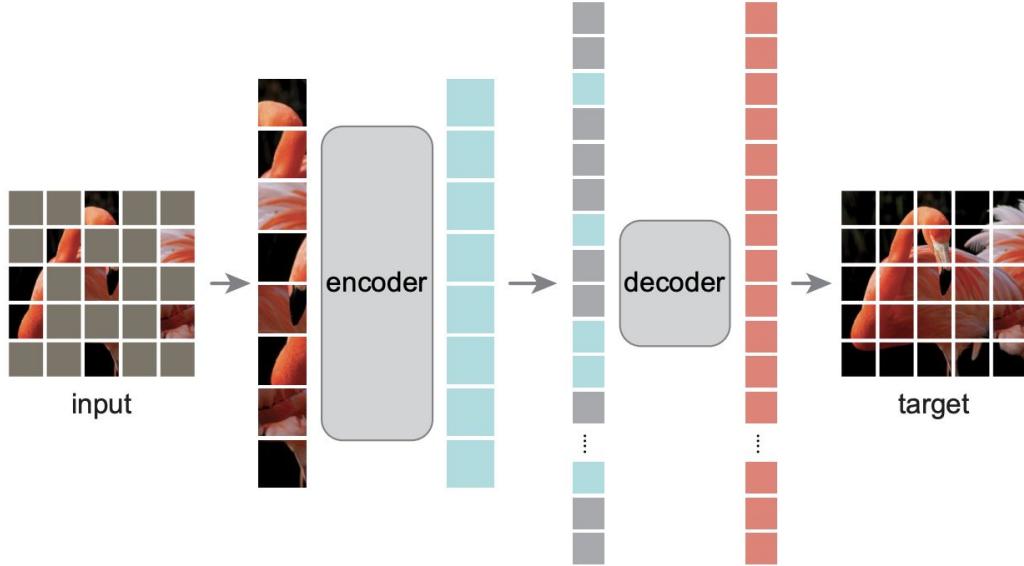


Figure 1. Our MAE architecture. During pre-training, a large random subset of image patches (*e.g.*, 75%) is masked out. The encoder is applied to the small subset of *visible patches*. Mask tokens are introduced *after* the encoder, and the full set of encoded patches and mask tokens is processed by a small decoder that reconstructs the original image in pixels. After pre-training, the decoder is discarded and the encoder is applied to uncorrupted images (full sets of patches) for recognition tasks.

¿Porque tanto % se enmascara?

Language:

Human-generated, Highly semantic and high density

Images: Natural Signals, heavy spatial redundancy

“Un patch faltante puede ser recuperado de los patches vecinos con poca comprensión semántica global de las partes, objetos y escenas”





original

mask 75%

mask 85%

mask 95%

MAE Encoder

- ViT
- Solo aplicado en los patches visibles (Eficiencia)
- Embeddings de los patches con una proyección lineal
- Positional embeddings added

MAE Decoder

- Tokens visibles del encoder
- Tokens enmascarados
- Positional embeddings a todos
- Una serie de bloques de transformers
- Decoder solo para pre-training (Tarea de reconstrucción)
- Arquitectura independiente del encoder (flexibilidad y eficiencia)
- 10% de procesamiento por token respecto al Encoder
- La última capa es lineal (dimensión transformer → cantidad de pixels/patch)
- Loss = MSE (Solo en patches enmascarados)

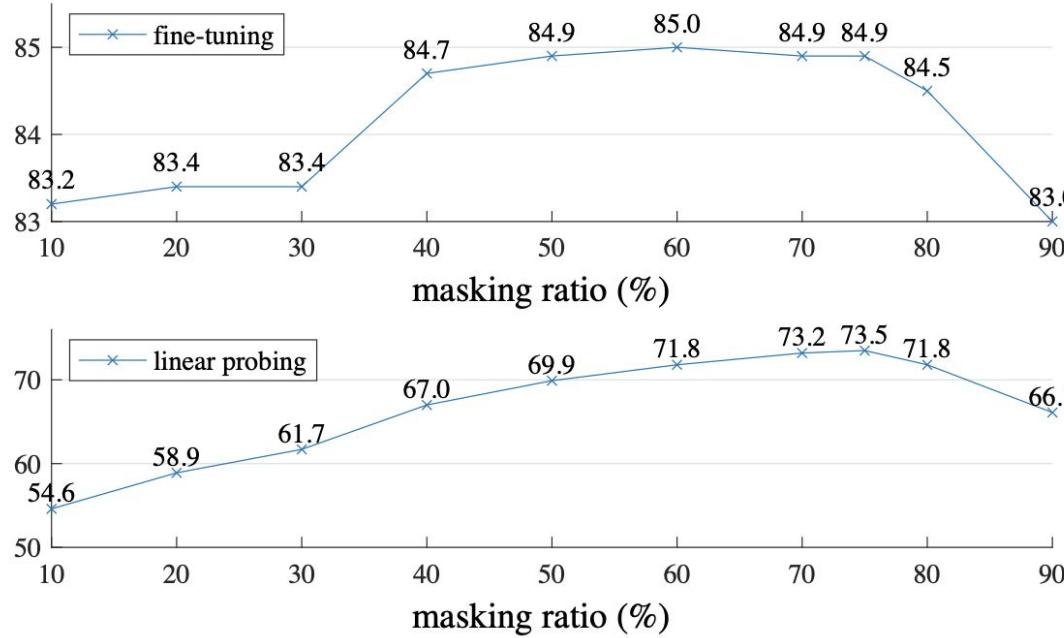


Figure 5. **Masking ratio.** A high masking ratio (75%) works well for both fine-tuning (top) and linear probing (bottom). The y-axes are ImageNet-1K validation accuracy (%) in all plots in this paper.

scratch, original [16]

scratch, our impl.

baseline MAE

76.5

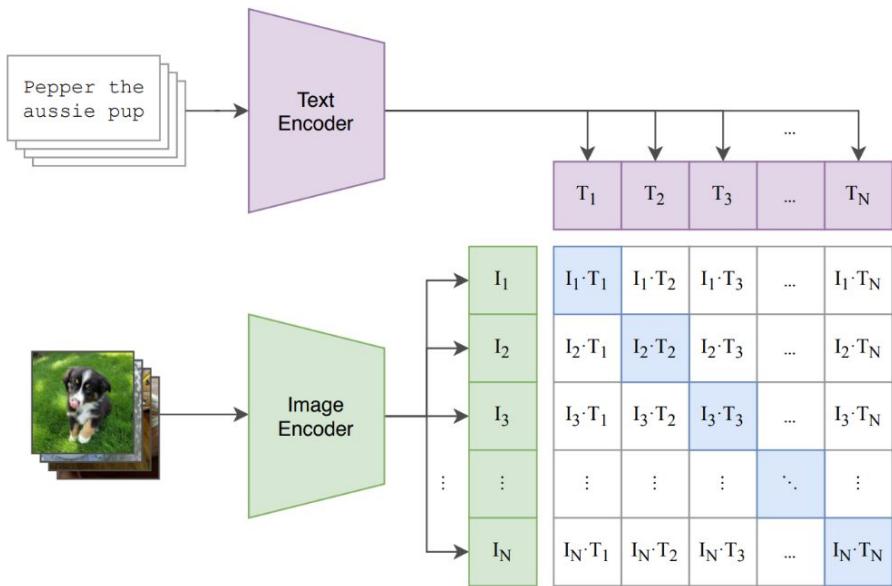
82.5

84.9

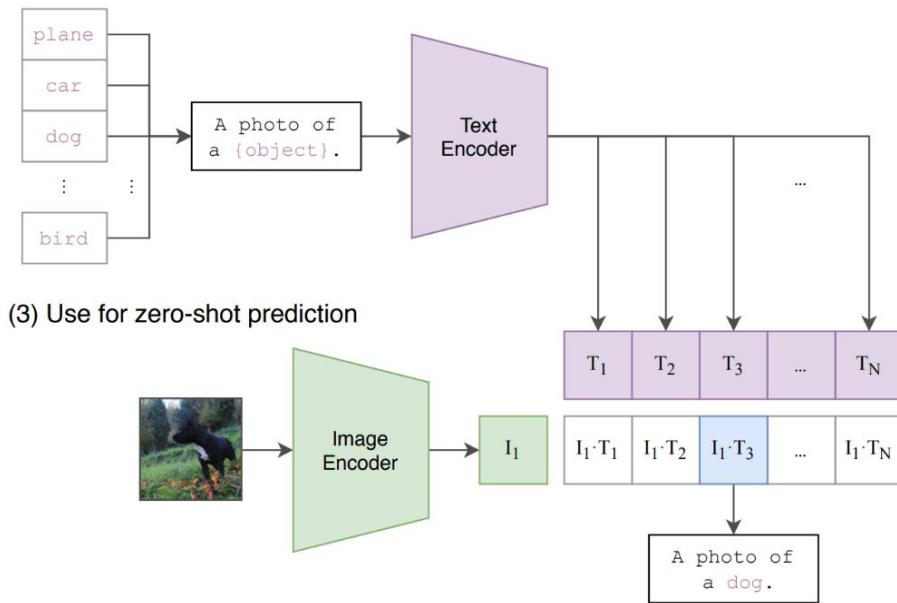
CLIP (Contrastive Language-Image Pre-Training)

[Learning Transferable Visual Models From Natural Language Supervision](#)
[repo openai](#)

(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

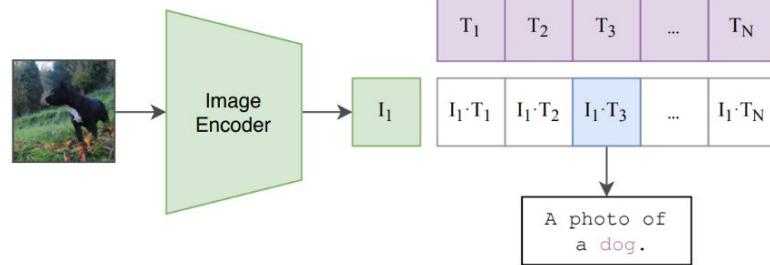


Figure 1. Summary of our approach. While standard image models jointly train an image feature extractor and a linear classifier to predict some label, CLIP jointly trains an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) training examples. At test time the learned text encoder synthesizes a zero-shot linear classifier by embedding the names or descriptions of the target dataset's classes.

Mini batch = 32K

Datasets:

Training: 400 millones de imágenes y texto (SOTA)

Benchmark:

- 30 CV datasets
- Tareas como: OCR, Reconocimiento de acciones (videos), Geo y otros
- Competitivo con fully supervised sin necesidad de entrenamiento específico
- Iguala la performance de Resnet50 en Imagenet (zero shot sin entrenar con 1.28 millones de imágenes)

Contrastive Loss ([Paper](#) - Yan Lecun et all)

$$D_W(\vec{X}_1, \vec{X}_2) = \|G_W(\vec{X}_1) - G_W(\vec{X}_2)\|_2 \quad (1)$$

$$\mathcal{L}(W) = \sum_{i=1}^P L(W, (Y, \vec{X}_1, \vec{X}_2)^i) \quad (2)$$

$$L(W, (Y, \vec{X}_1, \vec{X}_2)^i) = (1 - Y)L_S(D_W^i) + YL_D(D_W^i) \quad (3)$$

$$L(W, Y, \vec{X}_1, \vec{X}_2) = (1 - Y)\frac{1}{2}(D_W)^2 + (Y)\frac{1}{2}\{\max(0, m - D_W)\}^2 \quad (4)$$

Contrastive loss CLIP

- N (image, text) pairs (32K)
- Jointly training an image encoder and text encoder
- **Maximize the cosine similarity** of the image and text embeddings of the N **real pairs** in the batch
- **Minimizing the cosine similarity** of the embeddings of the $N^2 - N$ **incorrect pairings**.
- We optimize a **symmetric cross entropy loss** over these similarity scores.

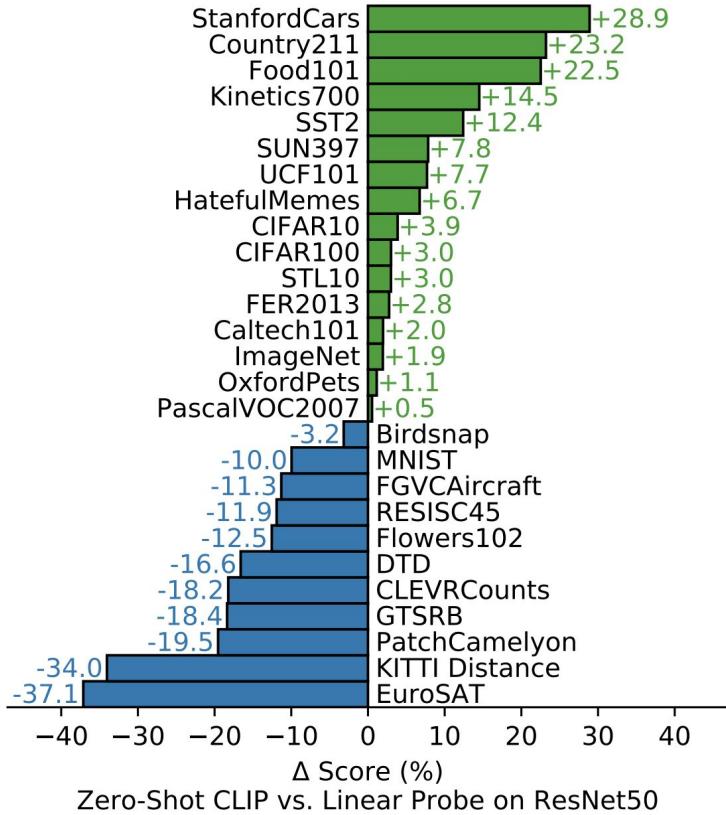


Figure 5. Zero-shot CLIP is competitive with a fully supervised baseline. Across a 27 dataset eval suite, a zero-shot CLIP classifier outperforms a fully supervised linear classifier fitted on ResNet-50 features on 16 datasets, including ImageNet.

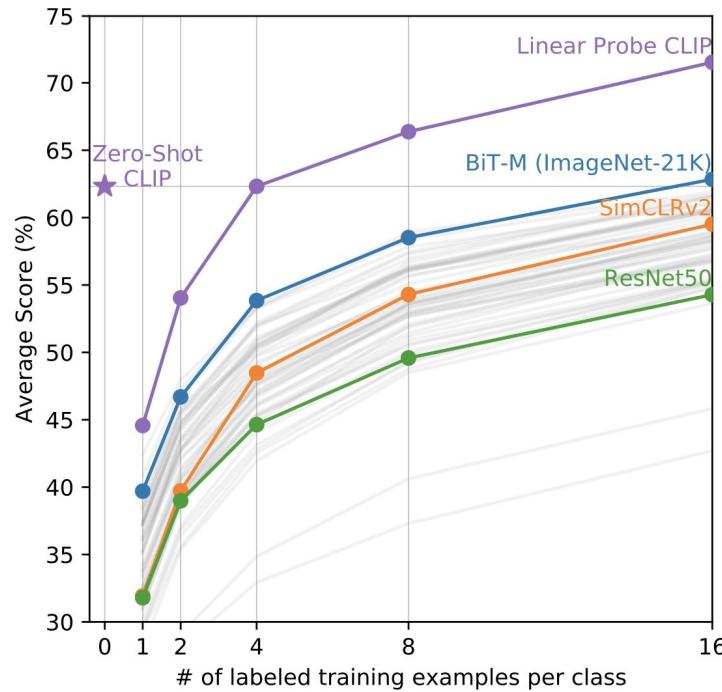


Figure 6. Zero-shot CLIP outperforms few-shot linear probes. Zero-shot CLIP matches the average performance of a 4-shot linear classifier trained on the same feature space and nearly matches the best results of a 16-shot linear classifier across publicly available models. For both BiT-M and SimCLRv2, the best performing model is highlighted. Light gray lines are other models in the eval suite. The 20 datasets with at least 16 examples per class were used in this analysis.

