





1

Transformers and foundation models

Malte Algren

Between Models and Reality





- A transformer takes in a set and returns a set
- Classification/regression: vector or scalar
- Add trainable latent space
 - Return that point at the end
 - Class token







- A transformer takes in a set and returns a set.
- Classification/regression: vector or scalar
- Add trainable latent space
 - Return that point at the end
 - Class token









- Full self attention also create a lot of junk
 - Where to send the junk?











- Full self attention also create a lot of junk
 - Where to send the junk?
 - Garbage collectors







- Full self attention also create a lot of junk
 - Where to send the junk?
- Adding trainable token that will be discarded afterwards











Foundation models How to make pretraining sound sexy

Malte Algren

Between Models and Reality



4



Google trends on "Foundation model" since 2020:

- Word2vec & GloVe (2013)
- GPT (2018)
- BERT/BEiT (2018)
- DALL-E (2021)

<u>A Survey on Self-supervised Learning</u> (2023)



Fig. 2: Google Scholar search results for "self-supervised learning". The vertical and horizontal axes denote the number of SSL publications and the year, respectively.





Google trends on "Foundation model" since 2020:

- Word2vec & GloVe (2013) ${}^{\bullet}$
- **GPT (2018)** ۲
- BERT/BEIT (2018) ullet
- DALL-E (2021) •

A Survey on Self-supervised Learning (2023)





Need to define some notation and language for foundation models

- 1. Backbone (encoder)
 - $g(x) = z \Longrightarrow f(z) = y'$
- Latent representation 2.
 - Representation: *z* ۲
- 3. Learning objective
 - Supervised: $\mathcal{L}(f(z), y)$
 - Self-supervised: $\mathcal{L}(f(z), x)$ •











Notation and definitions

Jan 3 2026

Subset of the lawsuits

Notation and definitions

ND California

Judge Chhabria

Judge Tigar

Judae Brever

Judge Alsup

Judge Lee

Judge Tigar

Judae Orrick

Zhang, Leovy)

Judae Lee

Bartz v. Anthropic

<u>Kadrey v. Meta; Chabon v. Meta</u>

In re OpenAl ChatGPT Litigation (Tremblay, Silverman, Chabon)

Judge Araceli Martínez-Olguír

Dubus v. NVIDIA Corp. (related)

O'Nan v. Databricks, Mosaic ML

Concord Music. v. Anthropic

<u>Sarah Andersen v. Stability Al</u>

In re Google Gen. Al Ltgn.

Makkai v. Databricks, Mosaic ML

Nazemian v. NVIDIA Corp.

Huckabee v. Meta; Farnsworth

Need to define some notation and language for foundation models

- 1. Backbone (encoder)
 - $g(x) = z \Longrightarrow f(z) = y'$
- 2. Latent representation
 - Representation: z •
- 3. Learning objective
 - Supervised: $\mathcal{L}(f(z), y)$
 - Self-supervised: $\mathcal{L}(f(z), x)$
- 4. Scale (The world biggest heist)
- 5. Downstream tasks $T_n(z)$ 11/06/2025





11





Pretraining method



- What does pretrain mean?
- Constructing a pretrain model to solve tasks $T_n(z)$
 - What is the object of a pretrained model?
 - 1. $T(z) \approx p(y)$
 - 2. $T(z) \approx p(x)$
 - 3. $T(z) \neq p(x) \lor p(y)$
 - 4. $T(z) \neq \neq \neq \neq \neq p(x) \lor p(y)$
 - Often T(z) is unknown but assume $T(z) \approx p(x)$









Pretraining method



z

 $\mathscr{L}(y,y')$

 $\blacktriangleright \mathscr{L}(x,x')$

 $\neg g(z)$ -

 $\begin{array}{c} \text{Decoder} \\ g(z) \end{array} \longrightarrow X'$

13





Pretraining method



- Constructing a pretrain model to solve tasks $T_n(z)$
 - 1. Supervised learning $T(z) \approx p(y)$
 - MSE, BCE, CE, etc.
 - Multi-task learning with σ
 - 2. Self-supervised $T(z) \approx p(x)$
 - 3. Semi-supervised $T(z) \approx p(x, y)$
 - Joint supervised and self-supervised
 - A. Final step: Finetuning









vector



 $\mathcal{L}(y,y')$

 $T_1(z)$

- If you have an idea about our downstream tasks ullet
 - Use x to model y access to y •
 - Multiple downstream task that are orthogonal
 - Board support in the phase space

Input sensor signals

Construct your representation z using $T_n(z)$ •



representation vectors

11/06/2025





- You don't have any y. You have access to x
 - No specific $T_n(z)$ in mind
 - Board support in the phase space
 - Always similar encoder/decoder setup
 - In physics we have a lot of record data but no labels







Self-supervised learning



- Constructing a task from samples X
- Augmentations is the name of the game

•
$$A(x) = x_a \Rightarrow \operatorname*{argmin}_{\theta} \mathcal{L}(f(x_a), x)$$

- 1. Masking
- 2. Cropping/cutout
- 3. Rotation
- 4. Other invariances
- Many of them are domain specific





Self-supervised learning



- Constructing a task from samples X •
- Augmentations is the name of the game •
 - Masking (general) ۲
 - Cropping (image) ${}^{\bullet}$



Generative

Joint embedding - CLR

- Masked autoencoder MAE 1. SimCLR
- **Diffusion MAE**

- - 2. BYOL
 - I-JEPA 3.

Compression or not?

- What do we want?
- Representation learning vs
- pretraining



Malte Algren



Masked autoencoder (wo/ compression)

- Apply masking to input and reconstruct it
- $A(x) = x_{input}, A^{-1}(x) = x_{target}$
- $\underset{\theta}{\operatorname{argmin}} \mathcal{L}(g_{\theta}(f_{\theta}(x_{input})), x_{target})$

Masking strategies and amount:



Without Compression



Without Compression







Masked autoencoder (wo/ compression)

- Apply masking to input and reconstruct it
- $A(x) = x_{input}, A^{-1}(x) = x_{target}$
- $\underset{\theta}{\operatorname{argmin}} \mathcal{L}(g_{\theta}(f_{\theta}(x_{input})), x_{target})$

Do you remember the issue with VAE?

Without Compression $\mathscr{L}(x,x')$ $\xleftarrow{}$ f(x) $\xrightarrow{}$ f(x) $\xrightarrow{}$ f(x) $\xrightarrow{}$ g(z) $\xrightarrow{}$ f(x) $\xrightarrow{}$ g(z) $\xrightarrow{}$ f(x) $\xrightarrow{}$ f(x) $\xrightarrow{}$ f(x) $\xrightarrow{}$ g(z) $\xrightarrow{}$ g(z) $\xrightarrow{}$ f(x) $\xrightarrow{}$ g(z) g(z) $\xrightarrow{}$ g(z) g(z)



11/06/2025





Diffusion Models as Masked Autoencoders

- Apply masking to input and reconstruct it
- $A(x) = x_{input}, A^{-1}(x) = x_{target}$
- $\underset{\theta}{\operatorname{argmin}} \mathcal{L}(g_{\theta}(f_{\theta}(x_{input})), x_{target})$
- Swap the decoder with a diffusion model
 - Same setup as MAE
 - Still no compression between $f_{\theta}(x_{input})$
 - But now compression can be introduced





Self-supervised learning: Generative

Malte Algren

Bottleneck Diffusion Models for Representation Learning

- $A(x) = x_{input}, A^{-1}(x) = x_{target}$
 - A is different representations of the image
- Diffusion autoencoder

LAS

- Continuous latent representation
- Allow for interpolation in the latent











Back to Masked AutoEncoders

- Drawbacks?
 - How do we ensure that the information we reconstruct is relevant?







Back to Masked autoencoders

- Drawbacks?
 - How do we ensure that the information we reconstruct is relevant?
- For image generation (VQ)-VAE are used to compress image
 - Then generate the image in the latent space









Compare features in the latent space

- 1. Augmentation in *x* to make x1 & x2
- 2. Encoder into latent space: z1 & z2
- 3. Similarity comparison
- Less redundant information in z1 & z2
- A lot of research in CL
 - SimCLR
 - BYOL

...







Compare features in the latent space

- 1. Augmentation in *x* to make x1 & x2
- 2. Encoder into latent space: z1 & z2
- 3. Similarity comparison
- Joint-Embedding Predictive Architecture







Compare features in the latent space

- 1. Augmentation in *x* to make x1 & x2
- 2. Encoder into latent space: z1 & z2
- 3. Similarity comparison
- Joint-Embedding Predictive Architecture
 - Extend to point clouds





Malte Algren





Generative

- 1. <u>Masked autoencoder</u> MAE
- 2. <u>Diffusion MAE</u>

Properties

- Loss/reconstruction in x
- Work very well
- Stable
- Have been sort of left by the by CS

Joint embedding

1. SimCLR

- 2. BYOL
- 3. I-JEPA

Properties

- Loss/reconstruction in z
- Very active when it comes to research
- Do also work well
- Can have model collapse

Malte Algren



Quick points on finetuning

Now you have a pretrained model

- How to finetune it for a task?
 - You have access to f(x) = z
- 1. Size of model
 - Can you even backpropagate through the model?
- 2. Domain shift
 - ssl in data, finetune on simulation

Thanks for listening! Questions?







