1.5em

# Basics of NLP and SSL
## Mathematical Foundations of ML 2023

Ido Ben-Shaul [1] [2], Shai Dekel [1]

[1]Tel-Aviv University, Israel

[2]eBay Research

June 10, 2023

# Transformers and Typical Tasks

Transformers and Typical Tasks

- Based on "Formal Algorithms for Transformers" [PH22]!

# Sequence Modelling (DTransformer)

- Given a vocabulary $V$, let $x_n \in V^*$ for $n \in [N_{\text{data}}]$ be a dataset of sequences "sampled" from dist. $P$. The goal is to learn an estimate $\hat{P}$ of the distribution $P(\mathbf{x})$.
- In practice, the distribution estimate is often decomposed via the chain rule as $\hat{P}(x) = \hat{P}_\theta(x[1]) \cdot \hat{P}_\theta(x[2] \mid x[1]) \cdots \hat{P}_\theta(x[\ell] \mid x[1 : \ell - 1])$, where $\theta$ consists of all neural network parameters to be learned.
- The goal is to learn a distribution over a single token $x[t]$ given its preceding tokens $x[1 : t - 1]$ as context.
- Examples include language modelling, RL policy distillation, or music generation.

# Sequence-to-sequence (seq2seq) prediction (EDTransformer)

- Given a vocabulary $V$ and an i.i.d. dataset of sequence pairs $(\boldsymbol{z}_n, \boldsymbol{x}_n) \sim P$, where $P$ is a distribution over $V^* \times V^*$, learn an estimate of the conditional distribution $P(\boldsymbol{x}|\boldsymbol{z})$.

- In practice, the conditional distribution estimate is often decomposed as $\hat{P}(\boldsymbol{x}|\boldsymbol{z}) = \hat{P}_\theta(x[1]\,|\,\boldsymbol{z}) \cdot \hat{P}_\theta(x[2]\,|\,x[1], \boldsymbol{z}) \cdots \hat{P}_\theta(x[\ell]\,|\,\boldsymbol{x}[1:\ell-1], \boldsymbol{z})$.

- Examples include translation ($\boldsymbol{z}$ = a sentence in English, $\boldsymbol{x}$ = the same sentence in German), question answering ($\boldsymbol{z}$ = question, $\boldsymbol{x}$ = the corresponding answer), text-to-speech ($\boldsymbol{z}$ = a piece of text, $\boldsymbol{x}$ = a voice recording of someone reading the text).

# Classification (ETransformer).

- Given a vocabulary $V$ and a set of classes $[N_C]$, let $(\boldsymbol{x}_n, c_n) \in V^* \times [N_C]$ for $n \in [N_{\text{data}}]$ be an i.i.d. dataset of sequence-class pairs sampled from $P(\boldsymbol{x}, c)$.
- The goal in classification is to learn an estimate of the conditional distribution $P(c|\boldsymbol{x})$.
- Examples include e.g. sentiment classification, spam filtering, toxicity classification.

# Tokenization: How Text is Represented

- In NLP, *tokenization* refers to how a piece of text such as "My grandma makes the best apple pie." is represented as a sequence of vocabulary elements (called *tokens*).
- **Character-level tokenization.** One possible choice is to let *V* be the English alphabet (plus punctuation).

## Tokenization: How Text is Represented

- In NLP, *tokenization* refers to how a piece of text such as "My grandma makes the best apple pie." is represented as a sequence of vocabulary elements (called *tokens*).
- **Character-level tokenization.** One possible choice is to let *V* be the English alphabet (plus punctuation).
- **Word-level tokenization.** *V* consists of all English words (plus punctuation).

# Tokenization: How Text is Represented

- In NLP, *tokenization* refers to how a piece of text such as "My grandma makes the best apple pie." is represented as a sequence of vocabulary elements (called *tokens*).
- **Character-level tokenization.** One possible choice is to let *V* be the English alphabet (plus punctuation).
- **Word-level tokenization.** *V* consists of all English words (plus punctuation).
- **Subword tokenization.** *V* is a set of commonly occurring word segments like 'cious', 'ing', 'pre'. Common words like 'is ' are often a separate token, and single characters are also included in *V* to ensure all words are expressible. E.g. [SHB15] used in GPT-2 [BMR$^+$20].

# Final vocabulary and text representation

- Given a choice of tokenization / vocabulary, each vocabulary element is assigned a unique index in $\{1, 2, \ldots, N_V - 3\}$.
- A number of special tokens are then added to the vocabulary. The number of special tokens varies, and here we will consider three:
  - mask_token $:= N_V - 2$, used in masked language modelling;
  - bos_token $:= N_V - 1$, used for representing the beginning of sequence;
  - eos_token $:= N_V$, used for representing the end of sequence.

  The complete vocabulary has $N_V = |V|$ elements.
- A piece of text is represented as a sequence of indices (called *token IDs*) corresponding to its (sub)words, preceded by bos_token and followed by eos_token.

# Notation

- Let $V$ denote a finite set, called a *vocabulary*, often identified with $[N_V] := \{1, ..., N_V\}$. This could be words or letters, but typically are sub-words, called tokens.
- Let $x \equiv x[1 : \ell] \equiv x[1]x[2] \ldots x[\ell] \in V^*$ be a sequence of tokens, e.g. a sentence or a paragraph or a document. Unlike in Python, we use arrays starting from 1, and $x[1 : \ell]$ includes $x[\ell]$.
- For a matrix $M \in \mathbb{R}^{d \times d'}$, we write $M[i, :] \in \mathbb{R}^{d'}$ for the $i$th row and $M[:, j] \in \mathbb{R}^d$ for the $j$-th column.
- The training data may naturally be a collection of (independent) articles, but even then, some may exceed the maximal context length $\ell_{\max}$ transformers can handle. In this case, an article is crudely broken into shorter chunks of length $\leq \ell_{\max}$.

Architectural Components

## Token embedding

The token embedding learns to represent each vocabulary element as a vector in $\mathbb{R}^{d_e}$.

---

**Algorithm** Token embedding.

**Input** : $v \in V \cong [N_V]$, a token ID.

**Output** : $e \in \mathbb{R}^{d_e}$, the vector representation of the token.

**Parameter:** $W_e \in \mathbb{R}^{d_e \times N_V}$, the token embedding matrix.

1 **return** $e = W_e[:, v]$

---

## Positional embedding

The positional embedding(PE) learns to represent a token's position in a sequence as in $\mathbb{R}^{d_e}$. E.g. the position of the 1st token in a sequence is represented by a (learned) vector $\boldsymbol{W_p}[:, 1]$, the position of the 2nd token is another (learned) vector $\boldsymbol{W_p}[:, 2]$, etc. Learned PE require that the input sequence length is at most some fixed number $\ell_{\max}$.

**Algorithm** Positional embedding.

**Input**     : $\ell \in [\ell_{\max}]$, position of a token in the sequence.
**Output**    : $\boldsymbol{e_p} \in \mathbb{R}^{d_e}$, the vector representation of the position.
**Parameter:** $\boldsymbol{W_p} \in \mathbb{R}^{d_e \times \ell_{\max}}$, the PE matrix.
2 **return** $\boldsymbol{e_p} = \boldsymbol{W_p}[:, \ell]$

The PE of a token is added to the token embedding to form a token's initial embedding. For the $t$-th token of a sequence $\boldsymbol{x}$, the embedding is

$$\boldsymbol{e} = \boldsymbol{W_e}[:, x[t]] + \boldsymbol{W_p}[:, t]. \tag{1}$$

Attention is the **main** part of transformers. It enables a neural network to make use of contextual information for predicting the current token.
On a high level:

- the token currently being predicted **(destination)** is mapped to a *query* vector $\boldsymbol{q} \in \mathbb{R}^{d_{\text{attn}}}$, and the tokens in the context **(source)** are mapped to *key* vectors $\boldsymbol{k}_t \in \mathbb{R}^{d_{\text{attn}}}$ and *value* vectors $\boldsymbol{v}_t \in \mathbb{R}^{d_{\text{value}}}$.
- The inner products $\boldsymbol{q}^\top \boldsymbol{k}_t$ are interpreted as the degree to which token **(src.)** $t$ is important for predicting the current token **(dst.)** $q$ – they are used to derive a distribution over the context tokens, which is then used to combine the value vectors.

# Attention



Multiplying x1 by the WQ weight matrix produces q1, the "query" vector associated with that word. We end up creating a "query", a "key", and a "value" projection of each word in the input sentence.
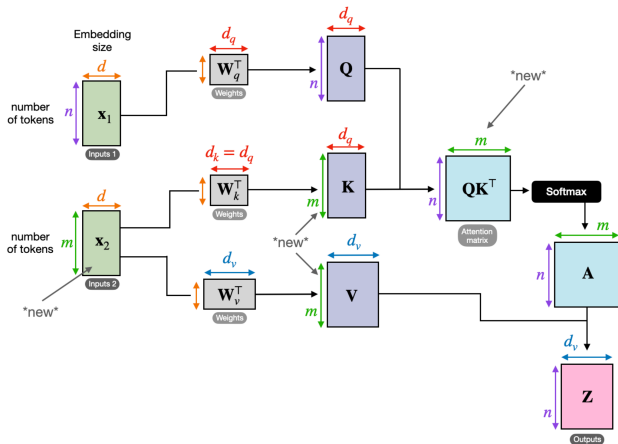
# Attention

Figure: Source

## Basic single-query attention.

**Algorithm** Basic Single-Query Attention

**Input**      : $e \in \mathbb{R}^{d_{in}}$, vector representation of the current token.

**Input**      : $e_t \in \mathbb{R}^{d_{in}}$, vector representations of context tokens $t \in [T]$.

**Output**   : $\tilde{v} \in \mathbb{R}^{d_{out}}$, vector representation of the token and context combined.

**Parameter:** $W_q, W_k \in \mathbb{R}^{d_{attn} \times d_{in}}$, the query and key linear projections.

**Parameter:** $W_v \in \mathbb{R}^{d_{out} \times d_{in}}$, the value linear projection.

3          $q \leftarrow W_q e$;

4          $\forall t: \; k_t \leftarrow W_k e_t$;

5          $\forall t: \; v_t \leftarrow W_v e_t$;

6          $\forall t: \; \alpha_t = \dfrac{\exp\left(q^\top k_t / \sqrt{d_{attn}}\right)}{\sum_{u \in [T]} \exp\left(q^\top k_u / \sqrt{d_{attn}}\right)}$;

7          **return** $\tilde{v} = \sum_{t=1}^{T} \alpha_t v_t$

# Additional Notation

It will be useful to define the softmax function for matrix arguments, as well as a Mask matrix:

- 
$$\text{softmax}(\boldsymbol{A})[t_z, t_x] := \frac{\exp A[t_z, t_x]}{\sum_t \exp A[t, t_x]}, \tag{2}$$

- 
$$\text{Mask}[t_z, t_x] = \begin{cases} 1 & \text{for bidirectional attention} \\ [[t_z \leq t_x]] & \text{for unidirectional att.} \end{cases} \tag{3}$$

**Algorithm** $\tilde{V} \leftarrow$ `Attention`$(X, Z | W_{qkv}, \text{Mask})$

/* Computes a single (masked) self- or cross-attention head. */

**Input:** $X \in \mathbb{R}^{d_x \times \ell_x}, Z \in \mathbb{R}^{d_z \times \ell_z}$, vector representations of primary and context sequence.

**Output:** $\tilde{V} \in \mathbb{R}^{d_{\text{out}} \times \ell_x}$, updated representations of tokens in $X$, folding in information from tokens in $Z$.

8  $W_{qkv}$ consisting of: $W_q \in \mathbb{R}^{d_{\text{attn}} \times d_x}$, $W_k \in \mathbb{R}^{d_{\text{attn}} \times d_z}$, $W_v \in \mathbb{R}^{d_{\text{out}} \times d_z}$, . Mask$\in \{0, 1\}^{\ell_z \times \ell_x}, \uparrow(3)$ ;

9     $Q \leftarrow W_q X$     $[\![Q\text{uery} \in \mathbb{R}^{d_{\text{attn}} \times \ell_x}]\!]$ ;

10     $K \leftarrow W_k Z$     $[\![K\text{ey} \in \mathbb{R}^{d_{\text{attn}} \times \ell_z}]\!]$ ;

11     $V \leftarrow W_v Z$ ;     $[\![V\text{alue} \in \mathbb{R}^{d_{\text{out}} \times \ell_z}]\!]$ ;

12     $S \leftarrow K^\top Q$ ; $\forall t_z, t_x$,     if $\neg \text{Mask}[t_z, t_x]$ then $S[t_z, t_x] \leftarrow -\infty$ ;

13     **return** $\tilde{V} = V \cdot \text{softmax}\left(S / \sqrt{d_{\text{attn}}}\right)$

# Attention Variants

- **Bidirectional / unmasked self-attention.** Given a sequence, attention to each token, treating all tokens in the sequence as the context. Algorithm 5, with $Z = X$ and no masking (Mask $\equiv 1$).
- **Unidirectional / masked self-attention.** Given a sequence, attention to each token, treating all preceding tokens (including itself) as the context. Future tokens are masked out, so this **causal auto-regressive** version can be used for online prediction. $Z = X$ and Mask $[t_z, t_x] := [\![t_z \leq t_x]\!]$. For this Mask, the output $\tilde{V}[:, 1 : t]$ only depends on $X[:, 1 : t]$, hence can be used to predict $X[:, t + 1]$.
- **Cross-attention.** Given two sequences (often in the context of a sequence-to-sequence task), attention to each token of the primary token sequence $X$, treating the second token sequence $Z$ as the context, with Mask $\equiv 1$. While the output $\tilde{V}$ and input sequences $X$ have the same length $\ell_x$, the context sequence $Z$ can have different length $\ell_z$.

---

**Algorithm** $\tilde{V} \leftarrow$ MHAttention($X, Z | W$, Mask)

---

/* Computes Multi-Head (masked) self- or cross- attention layer. */

**Input:** $X \in \mathbb{R}^{d_x \times \ell_x}, Z \in \mathbb{R}^{d_z \times \ell_z}$

**Output:** $\tilde{V} \in \mathbb{R}^{d_{out} \times \ell_x}$, updated representations of tokens in $X$, with information from tokens in $Z$.

14      **for** $h \in [H]$, $W_{qkv}^h$ *consisting of* **do**

15          $W_q^h \in \mathbb{R}^{d_{attn} \times d_x}, W_k^h \in \mathbb{R}^{d_{attn} \times d_z}, W_v^h \in \mathbb{R}^{d_{mid} \times d_z};$

16   $W_o \in \mathbb{R}^{d_{out} \times H d_{mid}}.$

17   $H$, number of attention heads, Mask$\in \{0, 1\}^{\ell_z \times \ell_x}$

18 **for** $h \in [H]$ **do**

19      $Y^h \leftarrow$ Attention($X, Z | W_{qkv}^h$, Mask)

20   $Y \leftarrow [Y^1; Y^2; \ldots; Y^H]$   **return** $\tilde{V} = W_o Y$

21 **return** $\tilde{V} = V \cdot \text{softmax} \left( S / \sqrt{d_{attn}} \right)$

---

## Layer Norm

Layer normalisation explicitly controls the mean and variance of individual neural network activations; the pseudocode is given in Algorithm 6.

---

**Algorithm** $\hat{e} \leftarrow \text{layer\_norm}(e|\gamma, \beta)$

---

/* Normalizes layer activations $e$. */

**Input** : $e \in \mathbb{R}^{d_e}$, neural network activations.

**Output** : $\widehat{e} \in \mathbb{R}^{d_e}$, normalized activations.

**Parameter:** $\gamma, \beta \in \mathbb{R}^{d_e}$, element-wise scale and offset.

2 $m \leftarrow \sum_{i=1}^{d_e} e[i]/d_e$;

3 $v \leftarrow \sum_{i=1}^{d_e} (e[i] - m)^2/d_e$;

4 **return** $\widehat{e} = \frac{e-m}{\sqrt{v}} \odot \gamma + \beta$, where $\odot$ denotes element-wise multiplication.

---

The unembedding learns to convert a vector representation of a token and its context into a distribution over the vocabulary elements

---

**Algorithm** Unembedding

---

**Input**     : $e \in \mathbb{R}^{d_e}$, a token encoding.
**Output**    : $p \in \Delta(V)$, a probability distribution over the vocabulary.
**Parameter:** $W_u \in \mathbb{R}^{N_V \times d_e}$, the unembedding matrix.
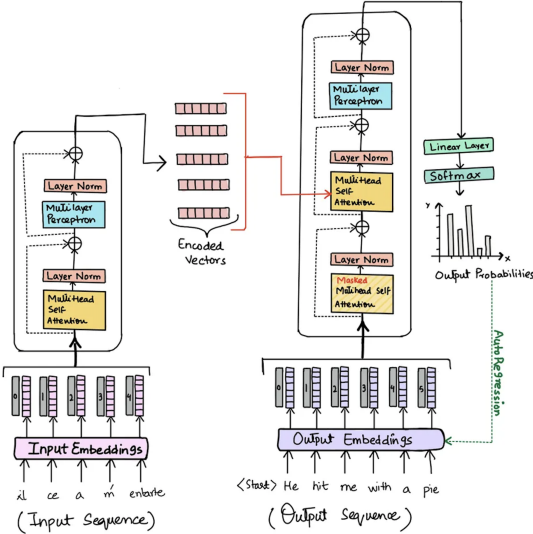25 **return** $p = \mathrm{softmax}(W_u e)$

---

Transformer Architectures

We will go over three example architectures:

- Encoder-Decoder Transformer (EDT) [VSP$^+$17]
- BERT (Encoder) [DCLT19]
- GPT (Decoder) [BMR$^+$20]

# EDT/ Sequence-to-sequence transformer

Intuition:

- First, the context sequence is encoded using bidirectional multi-head attention. The output of this 'encoder' part of the network is a vector representation of each context token, taking into account the entire context sequence.

- Second, the primary sequence is encoded. Each token in the primary sequence is allowed to use information from the encoded context sequence, as well as primary sequence tokens that precede it.

Figure: Left: T5 [RSR+19], Right: FLAN-T5 [CHL+22]

---

**Algorithm 8:** $P \leftarrow \text{EDTransformer}(z, x | \theta)$

---

/* Encoder-decoder transformer forward pass                                    */

**Input:** $z, x \in V^*$, two sequences of token IDs.

**Output:** $P \in (0, 1)^{N_V \times \text{length}(x)}$, where the $t$-th column of $P$ represents $\hat{P}_\theta(x[t+1] | x[1:t], z)$.

**Hyperparameters:** $\ell_{\max}, L_{\text{enc}}, L_{\text{dec}}, H, d_e, d_{\text{mlp}} \in \mathbb{N}$

/* Encode the context sequence:                                                */

1  $\ell_z \leftarrow \text{length}(z)$

2  for $t \in [\ell_z] : e_t \leftarrow W_e[:, z[t]] + W_p[:, t]$

3  $Z \leftarrow [e_1, e_2, \ldots e_{\ell_z}]$

4  **for** $l = 1, 2, \ldots L_{\text{enc}}$ **do**

5  $\quad$ $Z \leftarrow Z + \text{MHAttention}(Z | W_l^{\text{enc}}, \text{Mask} \equiv 1)$

6  $\quad$ for $t \in [\ell_z] : Z[:, t] \leftarrow \text{layer\_norm}(Z[:, t] | \gamma_l^1, \beta_l^1)$

7  $\quad$ $Z \leftarrow Z + W_{\text{mlp2}}^l \text{ReLU}(W_{\text{mlp1}}^l Z + b_{\text{mlp1}}^l \mathbf{1}^\top) + b_{\text{mlp2}}^l \mathbf{1}^\top$

8  $\quad$ for $t \in [\ell_z] : Z[:, t] \leftarrow \text{layer\_norm}(Z[:, t] | \gamma_l^2, \beta_l^2)$

9  **end**

# EDT - Architecture

```
/* Decode the primary sequence, conditioning on the context:          */
10  ℓ_x ← length(x)
11  for t ∈ [ℓ_x] :  e_t ← W_e[:, x[t]] + W_p[:, t]
12  X ← [e_1, e_2, ... e_{ℓ_x}]
13  for i = 1, 2, ..., L_dec do
14      X ← X + MHAttention(X ꞁ ᒽW_l^dec, Mask[t, t'] ≡ [[t ≤ t']])
15      for t ∈ [ℓ_x] :  X[:, t] ← layer_norm(X[:, t] ꞁ γ_l^3, β_l^3)
16      X ← X + MHAttention(X, Z ꞁ ᒽW_l^{e/d}, Mask ≡ 1)
17      for t ∈ [ℓ_x] :  X[:, t] ← layer_norm(X[:, t] ꞁ γ_l^4, β_l^4)
18      X ← X + W_{mlp4}^l ReLU(W_{mlp3}^l X + b_{mlp3}^l 1^⊤) + b_{mlp4}^l 1^⊤
19      for t ∈ [ℓ_x] :  X[:, t] ← layer_norm(X[:, t] ꞁ γ_l^5, β_l^5)
20  end
    /* Derive conditional probabilities and return:                    */
21  return P = softmax(W_u X)
```

---

**Algorithm 11:** $\hat{\theta} \leftarrow \text{EDTraining}(z_{1:N_{\text{data}}}, x_{1:N_{\text{data}}}, \theta)$

---

/* Training a seq2seq model                          */

**Input:** $\{(z_n, x_n)\}_{n=1}^{N_{\text{data}}}$, a dataset of sequence pairs.

**Input:** $\theta$, initial transformer parameters.

**Output:** $\hat{\theta}$, the trained parameters.

**Hyperparameters:** $N_{\text{epochs}} \in \mathbb{N}$, $\eta \in (0, \infty)$

1 **for** $i = 1, 2, \ldots, N_{\text{epochs}}$ **do**

2     **for** $n = 1, 2, \ldots N_{\text{data}}$ **do**

3         $\ell \leftarrow \text{length}(x_n)$

4         $P(\theta) \leftarrow \text{EDTransformer}(z_n, x_n | \theta)$

5         $\text{loss}(\theta) = -\sum_{t=1}^{\ell-1} \log P(\theta)[x_n[t+1], t]$

6         $\theta \leftarrow \theta - \eta \cdot \nabla \text{loss}(\theta)$

7     **end**

8 **end**

9 **return** $\hat{\theta} = \theta$

---

**Algorithm 15:** $\hat{x} \leftarrow \text{EDInference}(z, \hat{\theta})$

/* Using a trained seq2seq model for prediction.                                              */

**Input:** A seq2seq transformer and trained parameters $\hat{\theta}$ of the transformer.

**Input:** $z \in V^*$, input sequence, e.g. a sentence in English.

**Output:** $\hat{x} \in V^*$, output sequence, e.g. the sentence in German.

**Hyperparameters:** $\tau \in (0, \infty)$

1   $\hat{x} \leftarrow [\text{bos\_token}]$

2   $y \leftarrow 0$

3   **while** $y \neq eos\_token$ **do**

4      $P \leftarrow \text{EDTransformer}(z, \hat{x} | \hat{\theta})$

5      $p \leftarrow P[:, \text{length}(\hat{x})]$

6      sample a token $y$ from $q \propto p^{1/\tau}$

7      $\hat{x} \leftarrow [\hat{x}, y]$

8   **end**

9   **return** $\hat{x}$

# Encoder-only transformer: BERT

Intuition: BERT is a bidirectional transformer trained on the task of masked language modelling. Given a piece of text with some tokens masked out, the goal is to correctly recover the masked-out tokens. The original use of BERT was to learn generally useful text representations, which could then be adapted for various downstream NLP tasks. The masking is not performed via the Mask parameter but differently: During training each input token is replaced with probability $p_{mask}$ by a dummy token mask_token, and evaluation is based on the reconstruction probability of these knocked-out tokens.

# Encoder-only transformer: BERT, architecture

---

**Algorithm 9:** $P \leftarrow \text{ETransformer}(x|\theta)$

---

/* BERT, an encoder-only transformer, forward pass          */

**Input:** $x \in V^*$, a sequence of token IDs.

**Output:** $P \in (0,1)^{N_v \times \text{length}(x)}$, where each column of $P$ is a distribution over the vocabulary.

**Hyperparameters:** $\ell_{\max}, L, H, d_e, d_{\text{mlp}}, d_f \in \mathbb{N}$

**Parameters:** $\theta$ includes all of the following parameters:

    $W_e \in \mathbb{R}^{d_e \times N_v}$, $W_p \in \mathbb{R}^{d_e \times \ell_{\max}}$, the token and positional embedding matrices.

    For $l \in [L]$:

      | $\mathcal{W}_l$, multi-head attention parameters for layer $l$, see (4),

      | $\gamma_l^1, \beta_l^1, \gamma_l^2, \beta_l^2 \in \mathbb{R}^{d_e}$, two sets of layer-norm parameters,

      | $W_{\text{mlp1}}^l \in \mathbb{R}^{d_{\text{mlp}} \times d_e}$, $b_{\text{mlp1}}^l \in \mathbb{R}^{d_{\text{mlp}}}$, $W_{\text{mlp2}}^l \in \mathbb{R}^{d_e \times d_{\text{mlp}}}$, $b_{\text{mlp2}}^l \in \mathbb{R}^{d_e}$, MLP parameters.

    $W_f \in \mathbb{R}^{d_f \times d_e}, b_f \in \mathbb{R}^{d_f}, \gamma, \beta \in \mathbb{R}^{d_f}$, the final linear projection and layer-norm parameters.

    $W_u \in \mathbb{R}^{N_v \times d_e}$, the unembedding matrix.

1   $\ell \leftarrow \text{length}(x)$

2   for $t \in [\ell]$ : $e_t \leftarrow W_e[:, x[t]] + W_p[:, t]$

3   $X \leftarrow [e_1, e_2, \dots e_\ell]$

4   **for** $l = 1, 2, \dots, L$ **do**

5       $X \leftarrow X + \text{MHAttention}(X \,|\, \mathcal{W}_l, \text{Mask} \equiv 1)$

6       for $t \in [\ell]$ : $X[:, t] \leftarrow \text{layer\_norm}(X[:, t] \,|\, \gamma_l^1, \beta_l^1)$

7       $X \leftarrow X + W_{\text{mlp2}}^l \text{GELU}(W_{\text{mlp1}}^l X + b_{\text{mlp1}}^l \mathbf{1}^\intercal) + b_{\text{mlp2}}^l \mathbf{1}^\intercal$

8       for $t \in [\ell]$ : $X[:, t] \leftarrow \text{layer\_norm}(X[:, t] \,|\, \gamma_l^2, \beta_l^2)$

9   **end**

10   $X \leftarrow \text{GELU}(W_f X + b_f \mathbf{1}^\intercal)$

11   for $t \in [\ell]$ : $X[:, t] \leftarrow \text{layer\_norm}(X[:, t] \,|\, \gamma, \beta)$

12   **return** $P = \text{softmax}(W_u X)$

# Encoder-only: Training

---

**Algorithm 12:** $\hat{\theta} \leftarrow \texttt{ETraining}(x_{1:N_{\text{data}}}, \theta)$

---

/* Training by masked language
   modelling                                              */

**Input:** $\{x_n\}_{n=1}^{N_{\text{data}}}$, a dataset of sequences.

**Input:** $\theta$, initial encoder-only transformer
   parameters.

**Output:** $\hat{\theta}$, the trained parameters.

**Hyperparameters:** $N_{\text{epochs}} \in \mathbb{N}$, $\eta \in$
   $(0, \infty)$, $p_{\text{mask}} \in (0, 1)$

1   **for** $i = 1, 2, \ldots, N_{\text{epochs}}$ **do**

2     **for** $n = 1, 2, \ldots, N_{\text{data}}$ **do**

3       $\ell \leftarrow \text{length}(x_n)$

4       **for** $t = 1, 2, \ldots, \ell$ **do**

5         $\tilde{x}_n[t] \leftarrow \texttt{mask\_token}$ or $x_n[t]$
   randomly with probability
   $p_{\text{mask}}$ or $1 - p_{\text{mask}}$

6       **end**

7       $\tilde{T} \leftarrow \{t \in [\ell] : \tilde{x}_n[t] = \texttt{mask\_token}\}$

8       $P(\theta) \leftarrow \texttt{ETransformer}(\tilde{x}_n \,|\, \theta)$

9       $\text{loss}(\theta) = -\sum_{t \in \tilde{T}} \log P(\theta)[x_n[t], t]$

10      $\theta \leftarrow \theta - \eta \cdot \nabla \text{loss}(\theta)$

11    **end**

12 **end**

---

Figure: [DCLT19]

# Classification using Encoder Models



Figure: Source

# Encoder-only: Clustering in intermediate layers



Figure: NCC mismatch for Sequence-Classification datasets: **RTE**, **MRPC** and **CoLA**, using both vanilla (**solid**) and SVSL(**dashed**) losses. **Top:** train NCC mismatch, **Bottom:** test NCC mismatch. We show only a subset of the transformer blocks for clearness. The shaded pink background shows the TPT for the vanilla loss experiment, and the blue for the SVSL. The background is shaded purple at epochs when both experiments are in the TPT.

Figure: Left: Pretrained, Right: Finetuned [XQP⁺22]

# Decoder-only transformers: GPT

GPT-2 and GPT-3 are large language models developed by OpenAI.
They all have similar architectures and are trained by autoregressive
language modelling: Given an incomplete sentence or paragraph, the goal
is to predict the next token.

The main difference from BERT is that GPT use unidirectional attention
instead of bidirectional attention.

GPT-3 is identical except larger, and replaces dense attention in Line 6 by
sparse attention, i.e. each token only uses a subset of the full context.

**Algorithm 10:** $P \leftarrow \text{DTransformer}(x|\theta)$

/* GPT, a decoder-only transformer, forward pass                              */

**Input:** $x \in V^*$, a sequence of token IDs.

**Output:** $P \in (0,1)^{N_V \times \text{length}(x)}$, where the $t$-th column of $P$ represents $\hat{P}_\theta(x[t+1]|x[1:t])$.

**Hyperparameters:** $\ell_{\max}, L, H, d_e, d_{\text{mlp}} \in \mathbb{N}$

**Parameters:** $\theta$ includes all of the following parameters:

$\quad W_e \in \mathbb{R}^{d_e \times N_V}, W_p \in \mathbb{R}^{d_e \times \ell_{\max}}$, the token and positional embedding matrices.

$\quad$ For $l \in [L]$:

$\quad\quad |\; \mathcal{W}_l$, multi-head attention parameters for layer $l$, see (4),

$\quad\quad |\; \gamma_l^1, \beta_l^1, \gamma_l^2, \beta_l^2 \in \mathbb{R}^{d_e}$, two sets of layer-norm parameters,

$\quad\quad |\; W_{\text{mlp1}}^l \in \mathbb{R}^{d_{\text{mlp}} \times d_e}, b_{\text{mlp1}}^l \in \mathbb{R}^{d_{\text{mlp}}}, W_{\text{mlp2}}^l \in \mathbb{R}^{d_e \times d_{\text{mlp}}}, b_{\text{mlp2}}^l \in \mathbb{R}^{d_e}$,  MLP parameters.

$\quad \gamma, \beta \in \mathbb{R}^{d_e}$, final layer-norm parameters.

$\quad W_u \in \mathbb{R}^{N_V \times d_e}$, the unembedding matrix.

1 $\ell \leftarrow \text{length}(x)$

2 **for** $t \in [\ell]$ : $e_t \leftarrow W_e[:, x[t]] + W_p[:, t]$

3 $X \leftarrow [e_1, e_2, \ldots e_\ell]$

4 **for** $l = 1, 2, \ldots, L$ **do**

5 $\quad$ **for** $t \in [\ell]$ : $\tilde{X}[:, t] \leftarrow \text{layer\_norm}(X[:, t] \,|\, \gamma_l^1, \beta_l^1)$

6 $\quad X \leftarrow X + \text{MHAttention}(\tilde{X} \,|\, \mathcal{W}_l, \text{Mask}[t, t'] = [[t \le t']])$

7 $\quad$ **for** $t \in [\ell]$ : $\tilde{X}[:, t] \leftarrow \text{layer\_norm}(X[:, t] \,|\, \gamma_l^2, \beta_l^2)$

8 $\quad X \leftarrow X + W_{\text{mlp2}}^l \text{GELU}(W_{\text{mlp1}}^l \tilde{X} + b_{\text{mlp1}}^l \mathbf{1}^\top) + b_{\text{mlp2}}^l \mathbf{1}^\top$

9 **end**

10 **for** $t \in [\ell]$ : $X[:, t] \leftarrow \text{layer\_norm}(X[:, t] \,|\, \gamma, \beta)$

11 **return** $P = \text{softmax}(W_u X)$

# Decoder-only transformers: GPT - training

---

**Algorithm 13:** $\hat{\theta} \leftarrow \texttt{DTraining}(x_{1:N_{\text{data}}}, \theta)$

---

/* Training next token prediction */

**Input:** $\{x_n\}_{n=1}^{N_{\text{data}}}$, a dataset of sequences.

**Input:** $\theta$, initial decoder-only transformer parameters.

**Output:** $\hat{\theta}$, the trained parameters.

**Hyperparameters:** $N_{\text{epochs}} \in \mathbb{N}$, $\eta \in (0, \infty)$

1 **for** $i = 1, 2, \ldots, N_{\text{epochs}}$ **do**
2     **for** $n = 1, 2, \ldots N_{\text{data}}$ **do**
3         $\ell \leftarrow \text{length}(x_n)$
4         $P(\theta) \leftarrow \texttt{DTransformer}(x_n \,|\, \theta)$
5         $\text{loss}(\theta) = -\sum_{t=1}^{\ell-1} \log P(\theta)[x_n[t+1], t]$
6         $\theta \leftarrow \theta - \eta \cdot \nabla \text{loss}(\theta)$
7     **end**
8 **end**
9 **return** $\hat{\theta} = \theta$

---

# Decoder-only transformers: Inference

**Algorithm 14:** $y \leftarrow \mathtt{DInference}(x, \hat{\theta})$

/* Prompting a trained model and using it for prediction. */

**Input:** Trained transformer parameters $\hat{\theta}$.

**Input:** $x \in V^*$, a prompt.

**Output:** $y \in V^*$, the transformer's continuation of the prompt.

**Hyperparameters:** $\ell_{\text{gen}} \in \mathbb{N}, \ \tau \in (0, \infty)$

1   $\ell \leftarrow \text{length}(x)$

2   **for** $i = 1, 2, \ldots \ell_{\text{gen}}$ **do**

3      $P \leftarrow \mathtt{DTransformer}(x \,|\, \hat{\theta})$

4      $p \leftarrow P[:, \ell + i - 1]$

5      sample a token $y$ from $q \propto p^{1/\tau}$

6      $x \leftarrow [x, y]$

7   **end**

8   **return** $y = x[\ell + 1 : \ell + \ell_{\text{gen}}]$

# Vision Transformer



Figure: [DBK+21]

Figure: [RKH+21]

Figure: [RKH+21], Blogpost

# Zero/One/Few Shot Learning



Figure: [BMR+20]

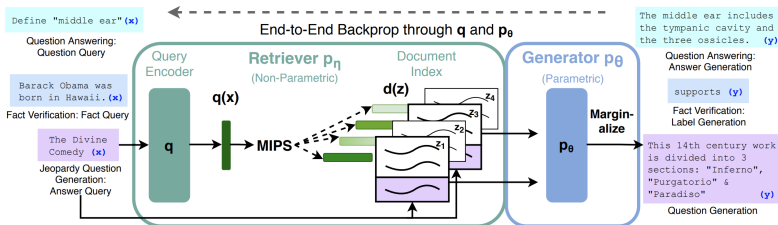# Chain-of-Thought prompting



Figure: [WWS+22]

Figure 1: Overview of our approach. We combine a pre-trained retriever (*Query Encoder + Document Index*) with a pre-trained seq2seq model (*Generator*) and fine-tune end-to-end. For query $x$, we use Maximum Inner Product Search (MIPS) to find the top-K documents $z_i$. For final prediction $y$, we treat $z$ as a latent variable and marginalize over seq2seq predictions given different documents.

GPT Assistant training pipeline

| Stage | Pretraining | Supervised Finetuning | Reward Modeling | Reinforcement Learning |
|-------|-------------|----------------------|-----------------|------------------------|
| Dataset | **Raw internet** text trillions of words low-quality, large quantity | **Demonstrations** Ideal Assistant responses, ~10-100K (prompt, response) written by contractors low quantity, high quality | **Comparisons** 100K –1M comparisons written by contractors low quantity, high quality | **Prompts** ~10K-100K prompts written by contractors low quantity, high quality |
| Algorithm | **Language modeling** predict the next token | **Language modeling** predict the next token | **Binary classification** predict rewards consistent w preferences | **Reinforcement Learning** generate tokens that maximize the reward |
| Model | Base model | SFT model *init from* | RM model *init from* | RL model *init from SFT use RM* |
| Notes | 1000s of GPUs months of training ex: GPT, LLaMA, PaLM **can deploy this model** | 1-100 GPUs days of training ex: Vicuna-13B **can deploy this model** | 1-100 GPUs days of training | 1-100 GPUs days of training ex: ChatGPT, Claude **can deploy this model** |

**You are a scientist who just invented a time machine. Where do you travel first?**

As a scientist, I would travel back in time to the year 1969 and watch the Apollo 11 mission land on Moon.

The Apollo 11 Mission was the first space flight that landed humans on the Moon. It was an incredibly important moment in history as it marked mankind's first step towards exploring other planets. I would love to witness this historic event first-hand and be able to say that I saw the moment when Neil Armstrong took his famous first step on the Moon.

**You are a scientist who just invented a time machine. Where do you travel first?**

I would travel to Ancient Egypt.

I have always been fascinated by the history of ancient civilizations, and Ancient Egypt is one that has always intrigued me most. Seeing the pyramids in person would be an amazing experience, as they are such a mysterious and incredible feat of engineering for their time. Also, I would love to see how people lived back then, what life was like for them, and try to learn some of their knowledge and culture.

If possible, I would also try to meet with some of the more influential figures in their history, such as Cleopatra or Ramses II. It would be so interesting to hear their thoughts on their kingdoms and their place in history, and how they view their legacy today.

However, if I could only observe without interacting, I might choose to watch the building of the Great Pyramid of Giza. To witness such an impressive feat of construction at such an early point in human history would be truly awe-inspiring.

Figure: LIMA: Less Is More for Alignment [ZLX$^+$23]

Self-Supervised Learning (SSL)

- How can we learn representations without labeling?

- How can we learn representations without labeling?
- **Trick:** Define a proxy task to learn representations
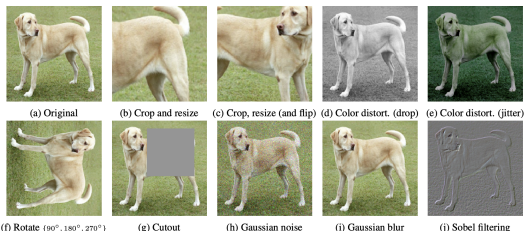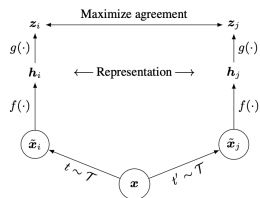- **Objective:** Learn meaningful representations without labels.

# SimCLR



Figure: SimCLR: [CKNH20]

$$L(f) = -\frac{1}{B} \sum_{i=1}^{B} \underbrace{(\text{sim}(Z_i, Z_i')/\tau)}_{\text{Invariance}} - \underbrace{\log \sum_{j \neq i} \exp(\text{sim}(Z_i, Z_j')/\tau)}_{\text{Regularization}},$$
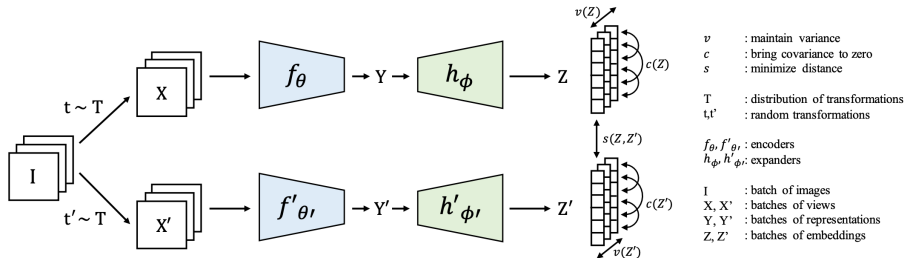
Figure: VICReg: [BPL22]

$$L(f) = \underbrace{\lambda s(Z, Z')}_{\text{Invariance}} + \underbrace{\mu[v(Z) + v(Z')] + \nu[c(Z) + c(Z')]}_{\text{Regularization}}, \qquad (4)$$

Understanding SSL

Figure: **SSL training induced semantic clustering.** UMAP of SSL representations in different hierarchies. **(top)** Augmentations of five different samples, each sample colored distinctly. **(middle)** Samples from five different classes within the standard CIFAR-100 dataset. **(bottom)** Samples from five different superclasses within the dataset. [BSSZG⁺23]

Figure: **SSL algorithms cluster the data with respect to semantic targets.**
The linear test accuracy rates, (left) non-normalized, (right) normalized by their
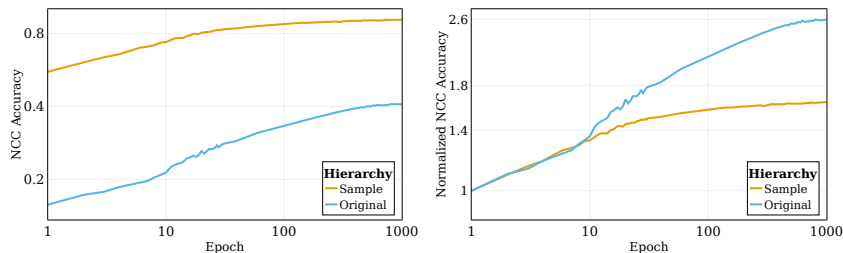values at initialization. All experiments are conducted on CIFAR-100 with VICReg
training

Figure: **SSL algorithms cluster the data with respect to semantic targets.**
The normalized NCC train accuracy, computed by dividing the accuracy values by
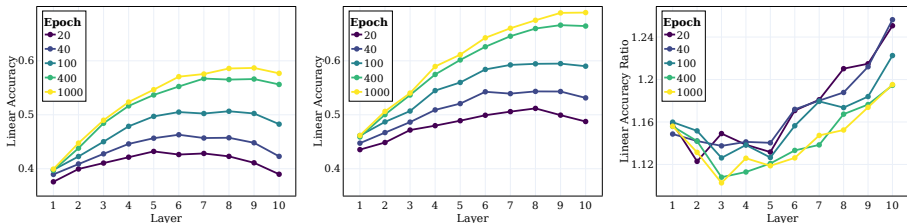their value at initialization.

Figure: **SSL efficiently learns semantic classes throughout intermediate layers.** The linear test accuracy of different layers of the model at various epochs **(left)** With respect to the 100 original classes. **(middle)** With respect to the 20 superclasses. **(right)** The ratio between the superclass and the original classes. All experiments are conducted on CIFAR-100 with VICReg training.

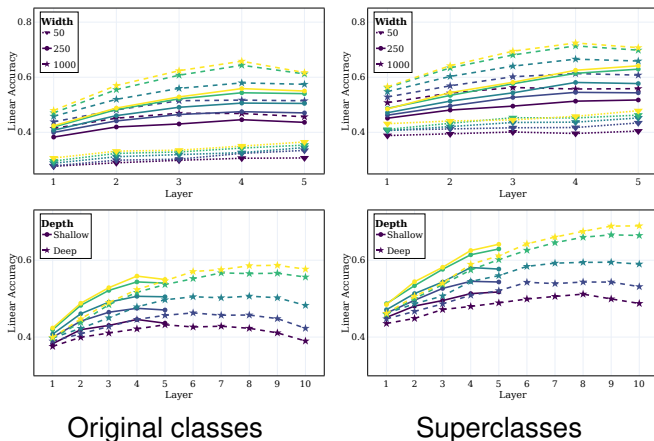Original classes                     Superclasses

Figure: **The influence of width and depth on learning semantic classes at intermediate layers. (top)** Linear test accuracy at different epochs for neural networks of varying widths. **(bottom)** Linear test accuracy of neural networks with different depths. **(left)** The performance is measured in relation to the original classes. **(right)** The performance with respect to the superclasses.
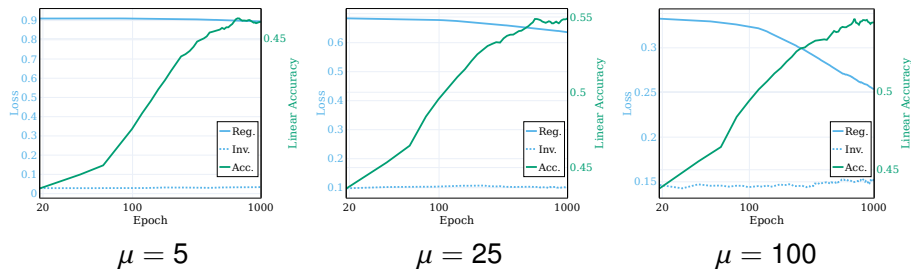
Figure: **The role of the regularization term in SSL training.** Each plot depicts the regularization and invariance losses, along with the linear test accuracy, throughout the training process of VICReg with $\mu = 5, 25, 100$ respectively.
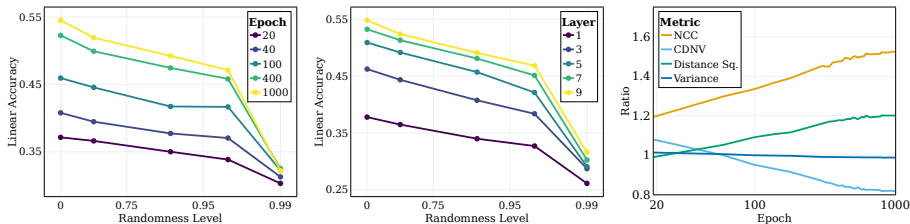
Figure: **SSL continuously learns semantic targets over random ones. (left)** The linear test accuracy for targets with varying levels of randomness from the last layers at different epochs. **(middle)** The linear test accuracy for targets with varying levels of randomness for the trained model. **(right)** The ratios between non-random and random targets for various clustering metrics. All experiments are conducted on CIFAR-100 with VICReg training.
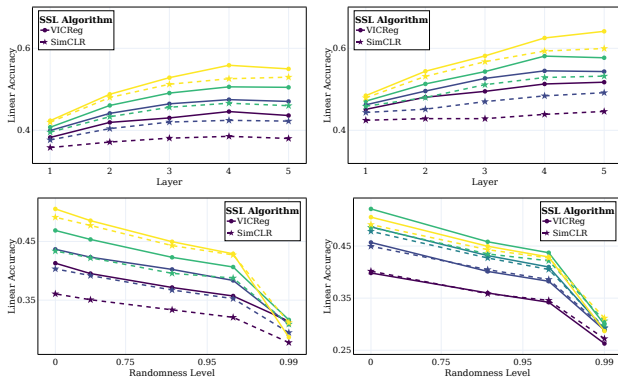
# Different Architectures



Figure: **SimCLR and VICReg have similar performance. (top)** Linear test accuracy in different training epochs, as a function of the intermediate layer, for original classes and superclasses, from left to right resp. **(bottom) (left)** Linear test accuracy in different training epochs (from dark to light) with respect to different randomness levels. **(right)** Linear test accuracy in different intermediate layers, at the end of training with respect to different randomness levels.
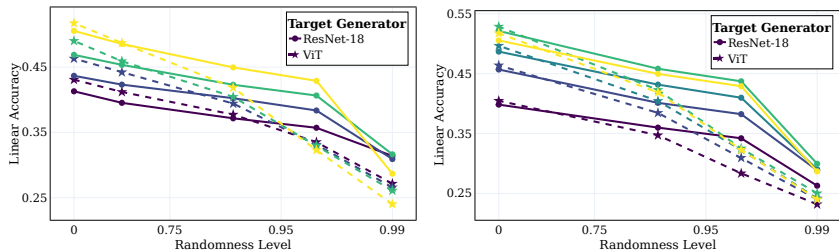
Figure: **The implicit bias of the backbone architecture on the learned representations. (left)** Linear test accuracy of an SSL-trained RES-5-250 network for extracting ResNet-18 and ViT random target functions with varying degrees of randomness (x-axis) at different epochs (color-coded from dark to bright). **(right)** Linear test accuracy of an SSL-trained RES-5-250 network for extracting ResNet-18 and ViT random target functions with varying degrees of randomness (x-axis) at different intermediate layers (color-coded from dark to bright).

# Bibliography

📄 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei.
Language models are few-shot learners.
In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.

📄 Adrien Bardes, Jean Ponce, and Yann LeCun.
VICReg: Variance-invariance-covariance regularization for self-supervised learning.
In *International Conference on Learning Representations*, 2022.

# Bibliography

📄 Ido Ben-Shaul, Ravid Shwartz-Ziv, Tomer Galanti, Shai Dekel, and Yann LeCun.
Reverse engineering self-supervised learning, 2023.

📄 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei.
Scaling Instruction-Finetuned Language Models.
*arXiv e-prints*, page arXiv:2210.11416, October 2022.

📄 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton.
A simple framework for contrastive learning of visual representations.
In *International Conference on Machine Learning (ICML)*, 2020.

📄 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby.
An image is worth 16x16 words: Transformers for image recognition at scale.
In *International Conference on Learning Representations*, 2021.

# Bibliography

📄 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova.
BERT: Pre-training of deep bidirectional transformers for language understanding.
In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.

📄 Mary Phuong and Marcus Hutter.
Formal Algorithms for Transformers.
*arXiv e-prints*, page arXiv:2207.09238, July 2022.

# Bibliography

📄 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
Learning Transferable Visual Models From Natural Language Supervision.
*arXiv e-prints*, page arXiv:2103.00020, February 2021.

📄 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu.
Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer.
*arXiv e-prints*, page arXiv:1910.10683, October 2019.

📄 Rico Sennrich, Barry Haddow, and Alexandra Birch.
Neural Machine Translation of Rare Words with Subword Units.
*arXiv e-prints*, page arXiv:1508.07909, August 2015.

# Bibliography

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin.
Attention is all you need.
In *Proceedings of the 30th International Conference on Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou.
Chain-of-Thought Prompting Elicits Reasoning in Large Language Models.
*arXiv e-prints*, page arXiv:2201.11903, January 2022.

# Bibliography

📄 Shuo Xie, Jiahao Qiu, Ankita Pasad, Li Du, Qing Qu, and Hongyuan Mei.
Hidden State Variability of Pretrained Language Models Can Guide Computation Reduction for Transfer Learning.
*arXiv e-prints*, page arXiv:2210.10041, October 2022.

📄 Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy.
LIMA: Less Is More for Alignment.
*arXiv e-prints*, page arXiv:2305.11206, May 2023.