Google DeepMind

# Language Modeling for Optimization

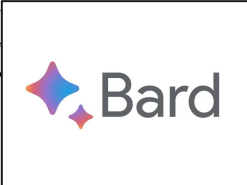Yutian Chen, Richard Song

13 September 2023

## Attention Is All You Need

## Language Models are Few-Shot Learners

Tom B. Brown*   Benjamin Mann*   Nick Ryder*   Melanie Subbiah*

Jared Kaplan†   Prafulla Dhariwal   Arvind Neelakantan   Pranav Shyam   Girish Sastry

Amanda Askell   Sandhini Agarwal   Ariel Herbert-

Daniel M. Zie

Eric Sigl

Jack Clark

dford

OpenAI

**Abstract**

reit*
rch
com

m

## Zero-Shot Text-to-Image Generation

Aditya Ramesh[1]   Mikhail Pavlov[1]   Gabriel Goh[1]   Scott Gray[1]
Chelsea Voss[1]   Alec Radford[1]   Mark Chen[1]   Ilya Sutskever[1]

**Abstract**

Text-to-image generation has traditionally fo-
cused on finding better
training on a fixed da
might involve compl
losses, or side informa
bels or segmentation n
ing. We describe a si
based on a transformer
els the text and image
data. With sufficient d
is competitive with pre
els when evaluated in

## AudioLM: a Language Modeling Approach to Audio Generation

Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin,
Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, Neil Zeghidour

*Abstract*—We introduce AudioLM, a framework for high-
quality audio generation with long-term consistency. AudioLM
maps the input audio to a sequence of discrete tokens and casts au-
relying on textual annotations. Yet, the acoustic diversity
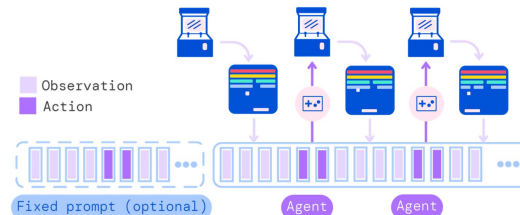and the quality remain limited: the model is trained on clean

## A Generalist Agent

Scott Reed*,†, Konrad Żołna*, Emilio Parisotto*, Sergio Gómez Colmenarejo†, Alexander Novikov,
Gabriel Barth-Maron, Mai Giménez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles,
Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals,
Mahyar Bordbar and Nando de Freitas†

*Equal contributions, †Equal senior contributions,  All authors are affiliated with DeepMind          reedscot@deepmind.com

# Transformers (Vaswani et al., 2017)

88K+ citations
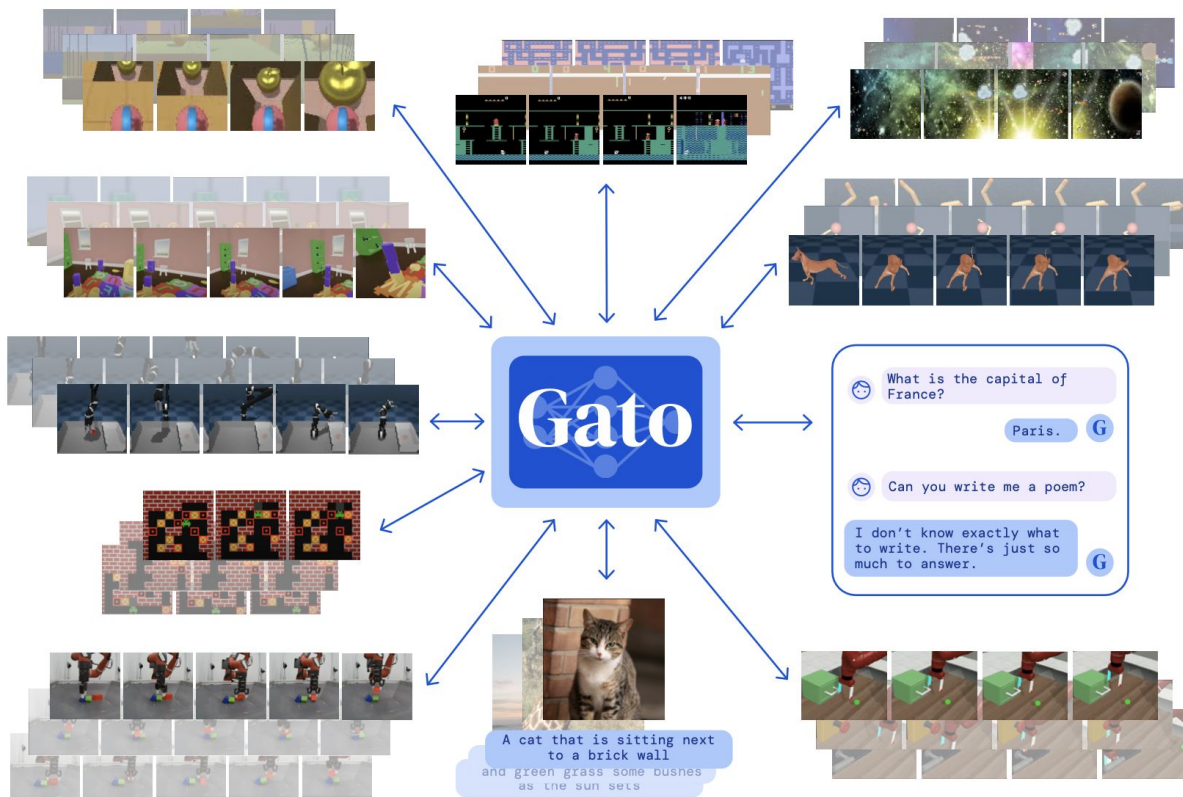
# In-context learning



outer loop

inner loop

**Learning via SGD during unsupervised pre-training**

| 1 | 5 + 8 = 13 |
| 2 | 7 + 2 = 9 |
| 3 | 1 + 0 = 1 |
| 4 | 3 + 4 = 7 |
| 5 | 5 + 9 = 14 |
| 6 | 9 + 8 = 17 |

↑
sequence #1

In-context learning

| 1 | gaot => goat |
| 2 | sakne => snake |
| 3 | brid => bird |
| 4 | fsih => fish |
| 5 | dcuk => duck |
| 6 | cmihp => chimp |

↑
sequence #2

In-context learning

| 1 | thanks => merci |
| 2 | hello => bonjour |
| 3 | mint => menthe |
| 4 | wall => mur |
| 5 | otter => loutre |
| 6 | bread => pain |

↑
sequence #3

In-context learning

GPT-3 (Brown et al., 2020)

# Unifying multi-modality as a single language model



Gato (Reed et al., 2022)

LLMs $\Bigg\{$
- NLP
- Speech
- Vision
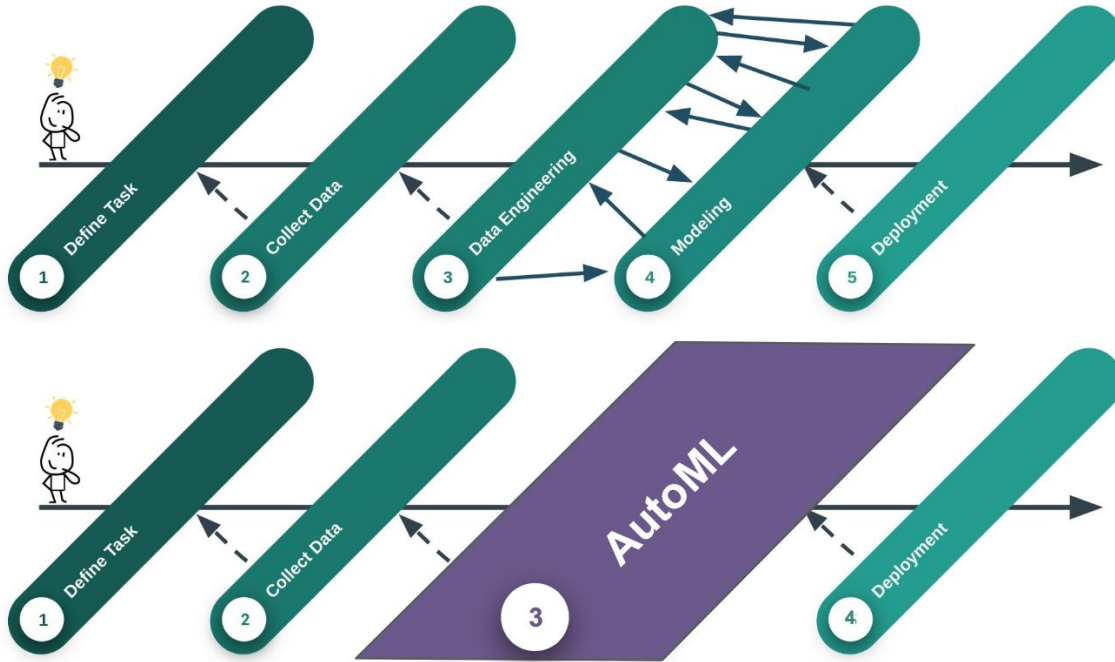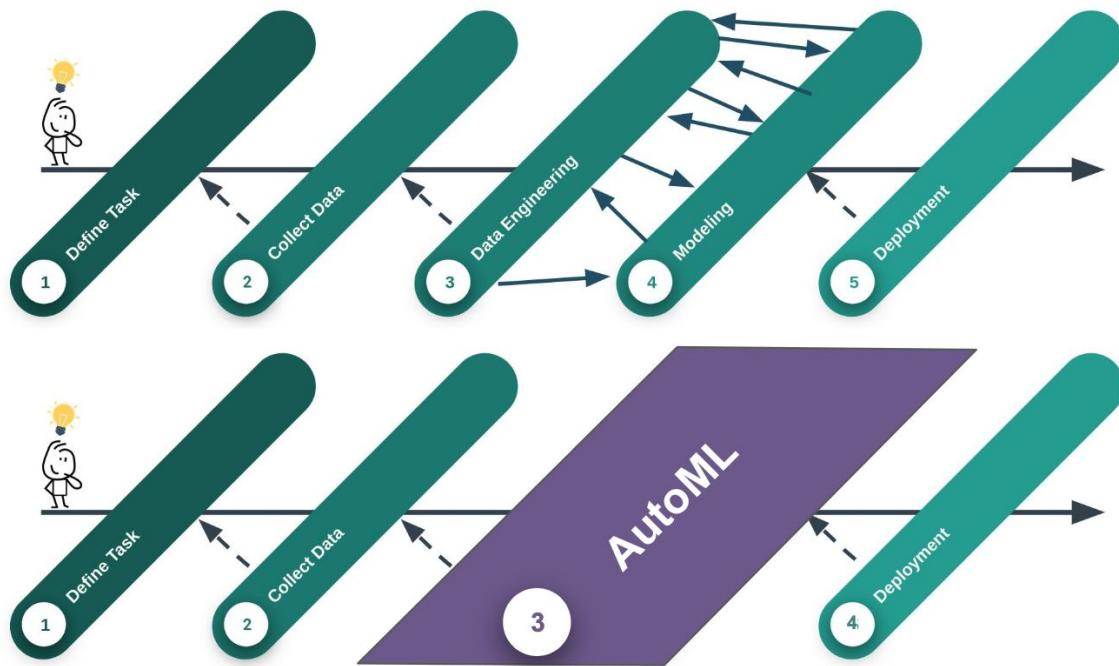- Games
- Controls
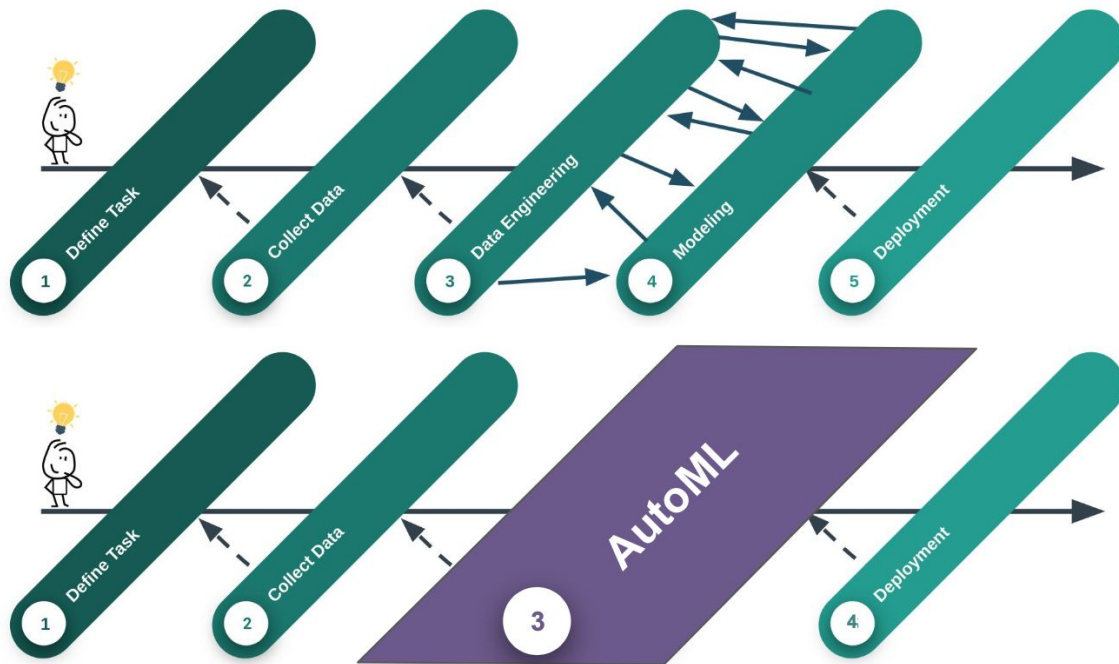- Robotics
- AI for Science
- …
- AutoML

# AutoML

# AutoML



**Traditional ML pipeline**

- *Data preprocess pipeline*
- *Feature engineering*
- *Select a model family*
- *Hyperparameters selection*
- *Model training*
- *Evaluation*

Figure source: https://www.automl.org/talks/

# AutoML



**AutoML**

- *Hyperparameter optimization (HPO)*
- *Neural Architecture Search (NAS)*
- *Meta–learning*
- *...*

Reference: https://www.automl.org/talks/

# Optimization in AutoML

$$\lambda^* = \arg\min_{\lambda \in \Lambda} \mathcal{L}(\lambda, \mathcal{D}_{\text{train}}, \mathcal{D}_{\text{valid}})$$

- Hyperparameter optimization (HPO)
    - λ: hyperparameter values
- Neural Architecture Search (NAS)
    - λ: network architecture
- Meta–learning
    - λ: optimization algorithm, model initialization, etc
- …

# Optimization in AutoML

$$\lambda^* = \arg\min_{\lambda \in \Lambda} \mathcal{L}(\lambda, \mathcal{D}_{\text{train}}, \mathcal{D}_{\text{valid}})$$

- Hyperparameter optimization (HPO)
  - λ: hyperparameter values
- Neural Architecture Search (NAS)
  - λ: network architecture
- Meta-learning
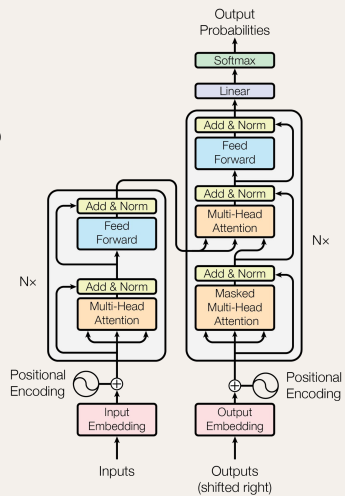  - λ: optimization algorithm, model initialization, etc
- …

*Gradient-free optimization*
*Black-box optimization (BBO)*

- *Bayesian optimization*
- *Reinforcement learning*
- *Evolutionary strategy*
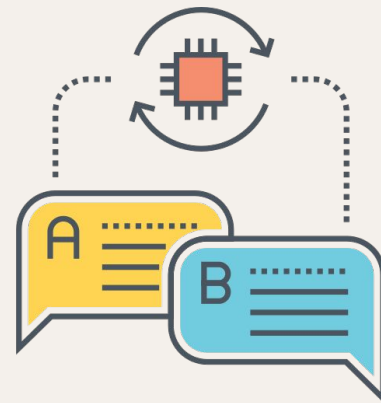- *Genetic algorithms*
- *…*

# What can LMs do for AutoML?

*Transformer architecture*

- *Powerful modeling capacity*
- *In-context learning*



*Go beyond BBO*

- *Leverage textual information*
- *More natural interface for AutoML*
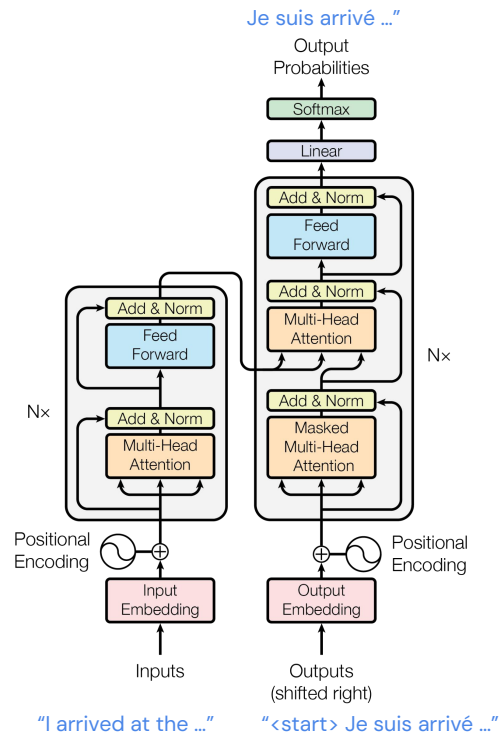
# Outline

Google DeepMind

- Basics of language modeling and transformers

- Transformers for optimization

- Large language models for optimization

- Vision and opening questions

**Google DeepMind**

# Outline

- **Basics of language modeling and transformers**

- Transformers for optimization

- Large language models for optimization

- Vision and opening questions

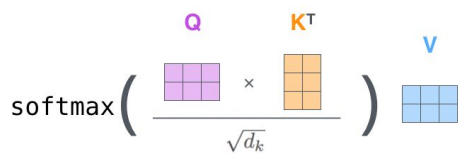# The Transformer: Attention is all you need
(Vaswani et al., 2017)

- Encoder-decoder architecture

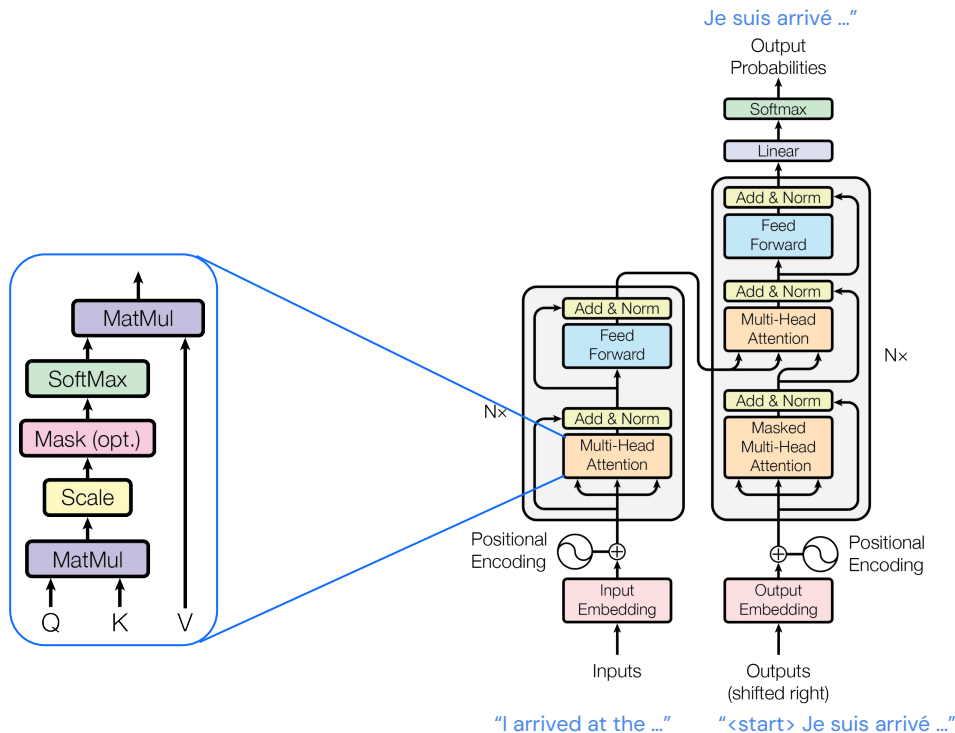# The Transformer: Attention is all you need
## (Vaswani et al., 2017)

- Encoder-decoder architecture

- Attention-base sequence model
  - Scaled Dot-Product Attention

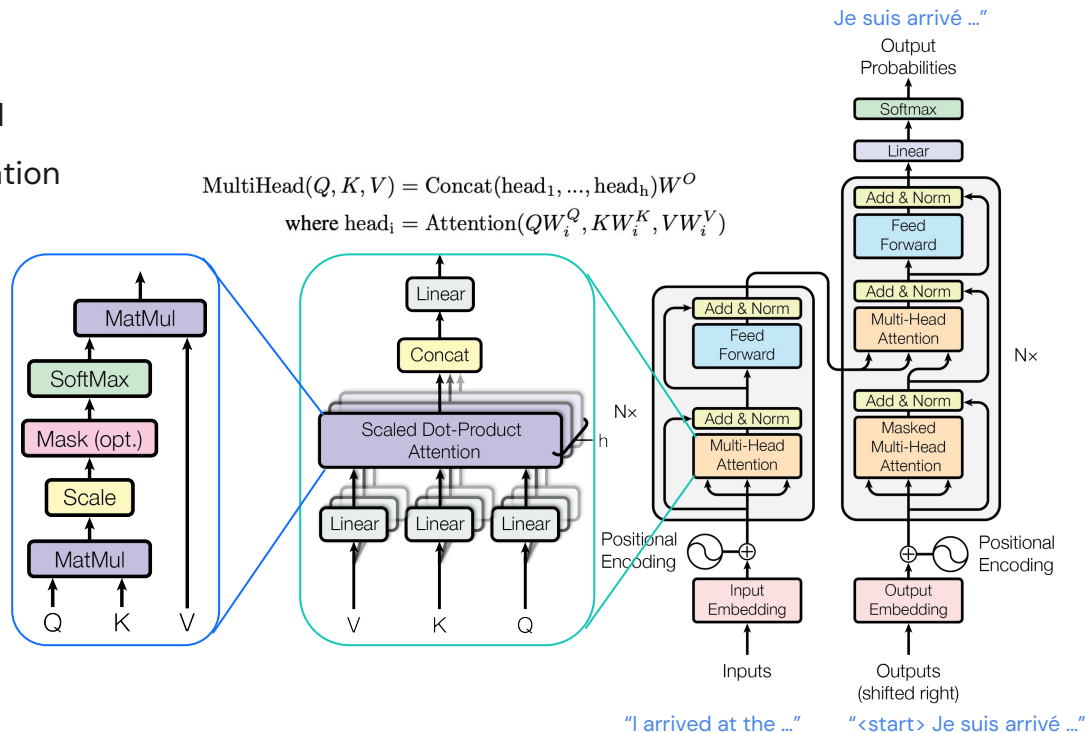$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$



(Figure source: jalammar.github.io/illustrated-transformer)

# The Transformer: Attention is all you need
## (Vaswani et al., 2017)

- Encoder-decoder architecture

- Attention-base sequence model
  - Scaled Dot-Product Attention
  - Multi-Head Attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$
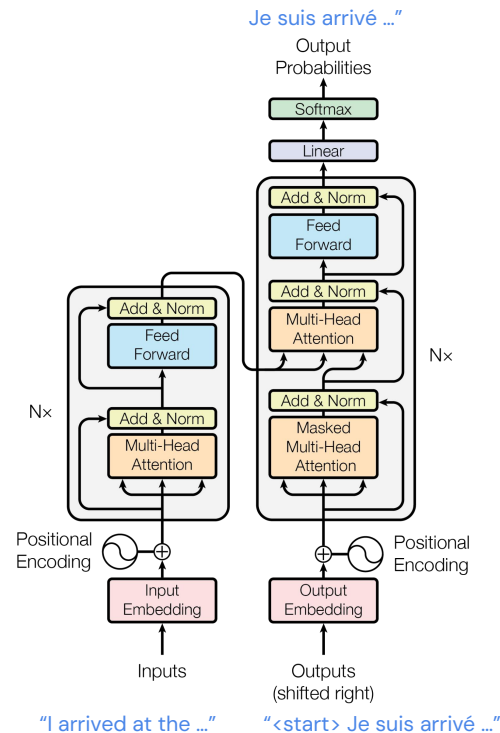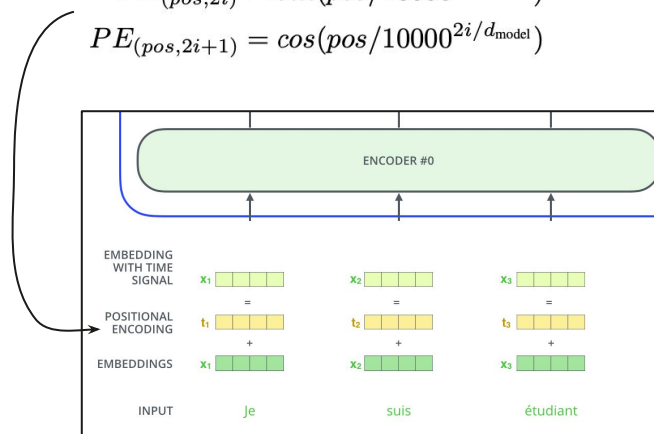$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

# The Transformer: Attention is all you need
(Vaswani et al., 2017)

- Encoder–decoder architecture

- Attention–base sequence model

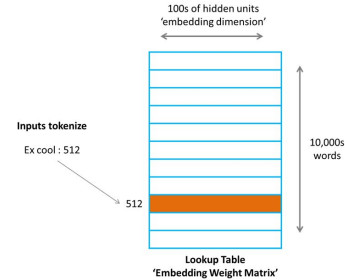- Positional encoding – representing the order of the sequence

$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{\mathrm{model}}})$$
$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{\mathrm{model}}})$$

# Transformers inputs

- **Text**

    - *Tokenizer*: sentence → token ids

        "I arrived at the ..." → 40, 5284, 379, 262, …

    - *Embedding table*: token ids → embedding vectors

        40, 5284, 379, 262, … →  Inputs = $[e_{40}, e_{5284}, e_{379}, e_{262}, …]$



100s of hidden units
'embedding dimension'

Inputs tokenize

Ex cool : 512

512

10,000s words

Lookup Table
'Embedding Weight Matrix'

# Transformers inputs
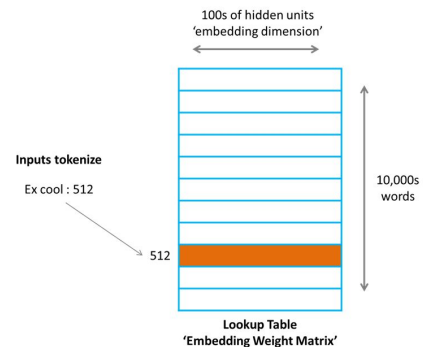
- **Text**
  - *Tokenizer*: sentence → token ids
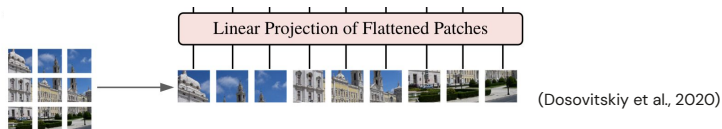
    "I arrived at the ..." → 40, 5284, 379, 262, ...

  - *Embedding table*: token ids → embedding vectors

    40, 5284, 379, 262, ... → Inputs = $[e_{40}, e_{5284}, e_{379}, e_{262}, ...]$



100s of hidden units
'embedding dimension'

Inputs tokenize

Ex cool : 512

512

10,000s words

Lookup Table
'Embedding Weight Matrix'

- **Images**



Linear Projection of Flattened Patches

(Dosovitskiy et al., 2020)

- **Audio**



Linear Projection
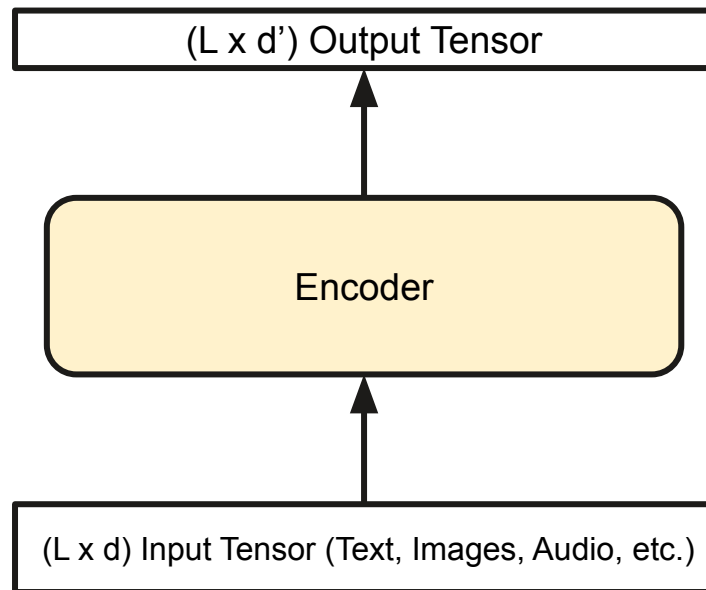
Input Spectrogram

Patch Split with Overlap

(Gong et al., 2021)

- **Discrete features**: One-hot encoding

  **Continuous features**: Normalization (linear, mu-law, Riemann distribution)

  (Chen et al., 2022; Reed et al., 2022; Muller et al., 2023)

# Transformers outputs: representation learning

**Use raw output (L x d') in any way**
- Directly into downstream models
- Take single slice (e.g. embeddings)

(L x d') Output Tensor

Encoder

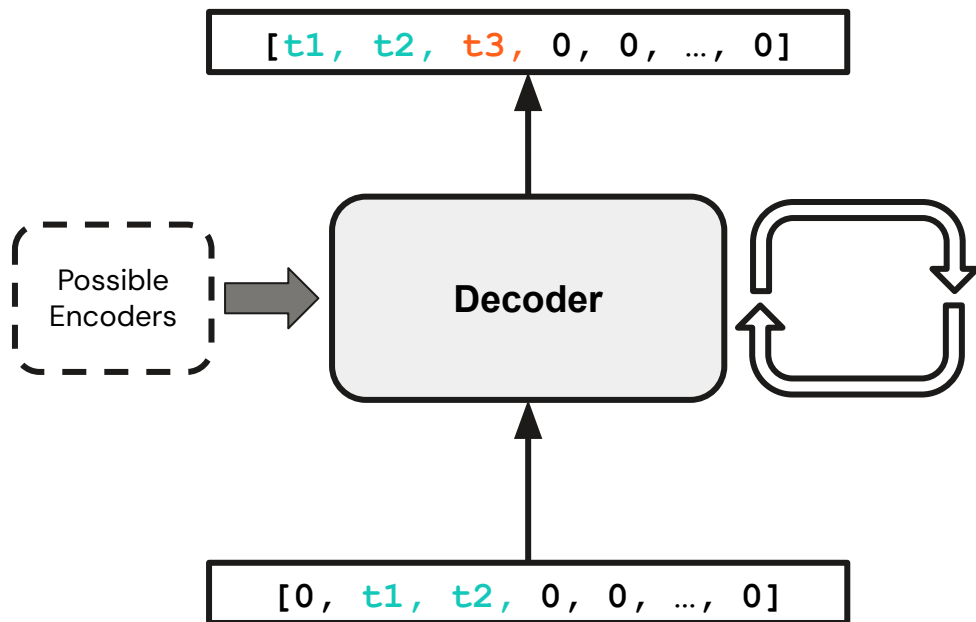(L x d) Input Tensor (Text, Images, Audio, etc.)

# Transformers outputs: answer generation

**Multi-step decoding:** 1 forward pass = **1 new token**

**Decoding Methods**

- Ancestral Sampling
  - Temperature sampling
  - Top-K Sampling
  - Nucleus Sampling
  - …

- Max likelihood sequence
  - Beam Search

`[t1, t2, t3, 0, 0, …, 0]`

Possible Encoders → **Decoder**
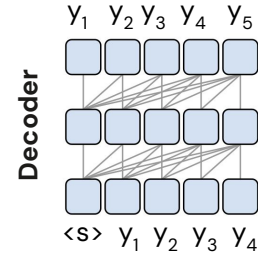
`[0, t1, t2, 0, 0, …, 0]`

# Transformers pre-training

- **Causal language modeling** (e.g. GPT, PaLM)

$$\mathcal{L} = \sum_i \log P\left(y_i | \text{Decoder}(\mathbf{y}_{1:i-1})\right)$$

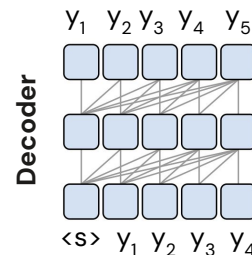   - E.g. log P("Thank you for inviting me to your party last week")

# Transformers pre-training



- **Causal language modeling** (e.g. GPT, PaLM)

  $$\mathcal{L} = \sum_i \log P\left(y_i | \text{Decoder}(\mathbf{y}_{1:i-1})\right)$$

  - E.g. log P("Thank you for inviting me to your party last week")

- **Masked language modeling** (e.g. BERT (Kenton and Toutanova, 2019))

  $$\mathcal{L} = \sum_{i \in \text{masked set}} \log P(y_i | \text{Encoder}(\mathbf{y}_{\text{masked}}))$$

  - E.g. log P("* * for inviting * * * * last *" | "Thank you * * me to your party * week")

# Transformers pre-training

- **Causal language modeling** (e.g. GPT, PaLM)

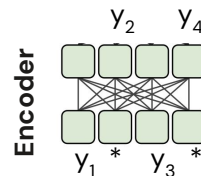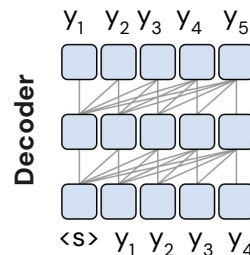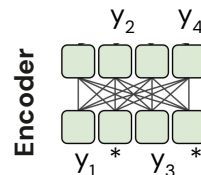  $$\mathcal{L} = \sum_i \log P\left(y_i | \text{Decoder}(\mathbf{y}_{1:i-1})\right)$$

  - E.g. log P("Thank you for inviting me to your party last week")

- **Masked language modeling** (e.g. BERT (Kenton and Toutanova, 2019))

  $$\mathcal{L} = \sum_{i \in \text{masked set}} \log P(y_i | \text{Encoder}(\mathbf{y}_{\text{masked}}))$$
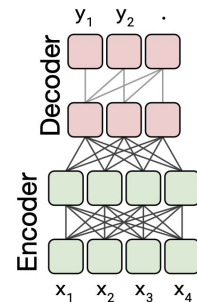
  - E.g. log P("* * for inviting * * * * last *" | "Thank you * * me to your party * week")

- **Hybrid objective** (e.g. T5 (Raffel et al., 2020))

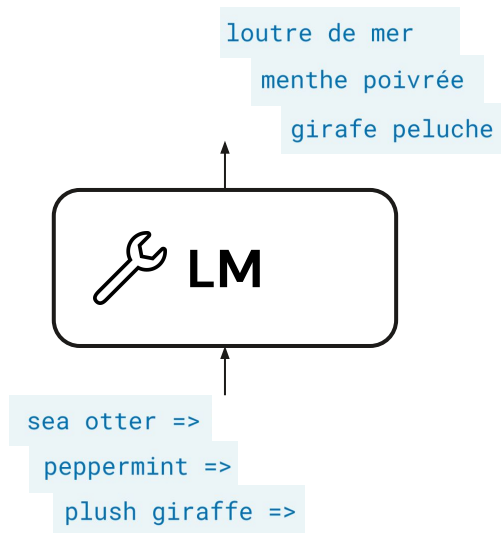  $$\mathcal{L} = \sum_i \log P\left(y_i | \text{Decoder}(\text{Encoder}(\mathbf{x}), \mathbf{y}_{1:i-1})\right)$$

  - E.g. log P("<X> for inviting <Y> last" | "Thank you <X> me to your party <Y> week")

# Downstream applications
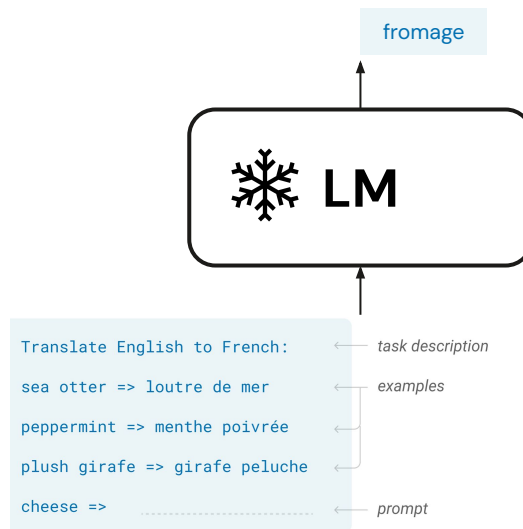
- **Fine-tuning**
    - Fine-tuning parameters on new data
    - Supervised FT
    - Reinforcement learning FT (RLHF)

```
loutre de mer
  menthe poivrée
    girafe peluche
```

🔧 **LM**

```
sea otter =>
  peppermint =>
    plush giraffe =>
```

- **Prompting**
    - Give task description and new data as *context* in the input

```
fromage
```

❄️ **LM**

```
Translate English to French:          ← task description
sea otter => loutre de mer            ← examples
peppermint => menthe poivrée          ←
plush girafe => girafe peluche        ←
cheese =>                             ← prompt
```

# Outline

Google DeepMind

# Optimization in AutoML

$$\lambda^* = \arg\min_{\lambda \in \Lambda} \mathcal{L}(\lambda, \mathcal{D}_{\text{train}}, \mathcal{D}_{\text{valid}})$$

- Hyperparameter optimization
  - $\lambda$: hyperparameter values
- Neural Architecture Search (NAS)
  - $\lambda$: network architecture
- Meta-learning
  - $\lambda$: optimization algorithm, model initialization, etc
- …

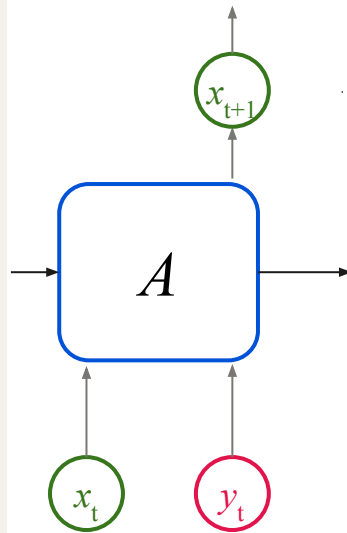*Gradient-free optimization*
*Black-box optimization (BBO)*

- *Bayesian optimization*
- *Reinforcement learning*
- *Evolutionary strategy*
- *Genetic algorithms*
- *…*

# Iterative optimization as sequence modeling

$$x^* = \arg\min_{x \in \mathcal{X}} f(x)$$

- An iterative optimizer

$$\mathcal{A} : x_t, f(x_t), h_t \to x_{t+1}$$

# Iterative optimization as sequence modeling

$$x^* = \underset{x \in \mathcal{X}}{\arg\min} \ f(x)$$

- An iterative optimizer
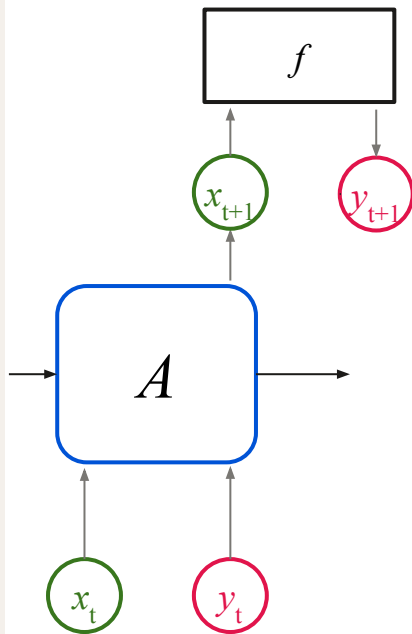
$$\mathcal{A} : x_t, f(x_t), h_t \rightarrow x_{t+1}$$

# Iterative optimization as sequence modeling

$$x^* = \arg\min_{x \in \mathcal{X}} f(x)$$

- An iterative optimizer

$$\mathcal{A} : x_t, f(x_t), h_t \to x_{t+1}$$
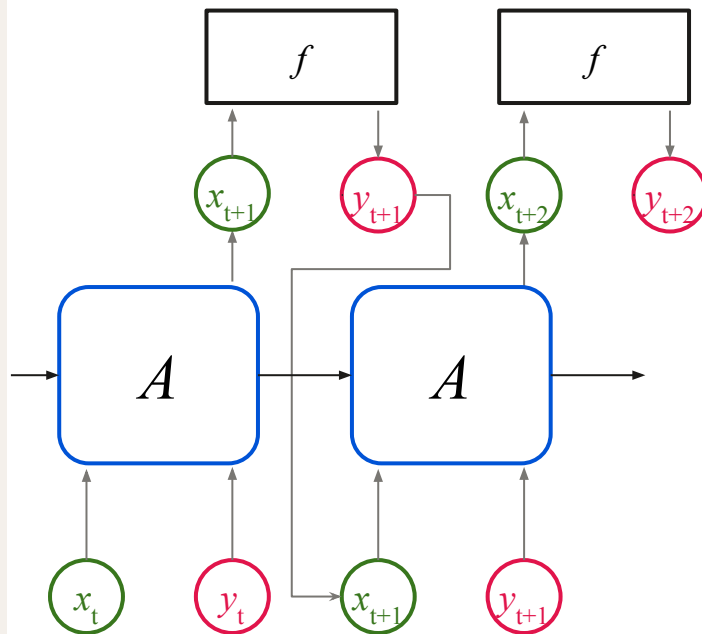
# Iterative optimization as sequence modeling

$$x^* = \arg\min_{x \in \mathcal{X}} f(x)$$

- An iterative optimizer

$$\mathcal{A} : x_t, f(x_t), h_t \rightarrow x_{t+1}$$

# Iterative optimization as sequence modeling

$$x^* = \arg\min_{x \in \mathcal{X}} f(x)$$

- An iterative optimizer

$$\mathcal{A} : x_t, f(x_t), h_t \rightarrow x_{t+1}$$

- Sequence–input optimizer
  (Chen et al., 2017)

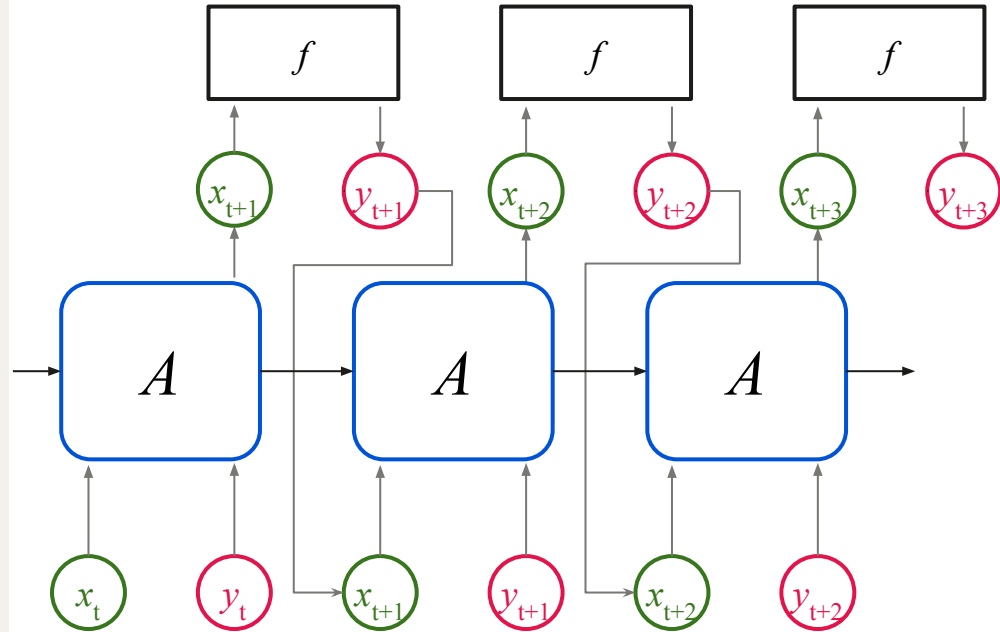$$\mathcal{A} : \mathbf{x}_{1:t}, \mathbf{y}_{1:t} \rightarrow x_{t+1}$$
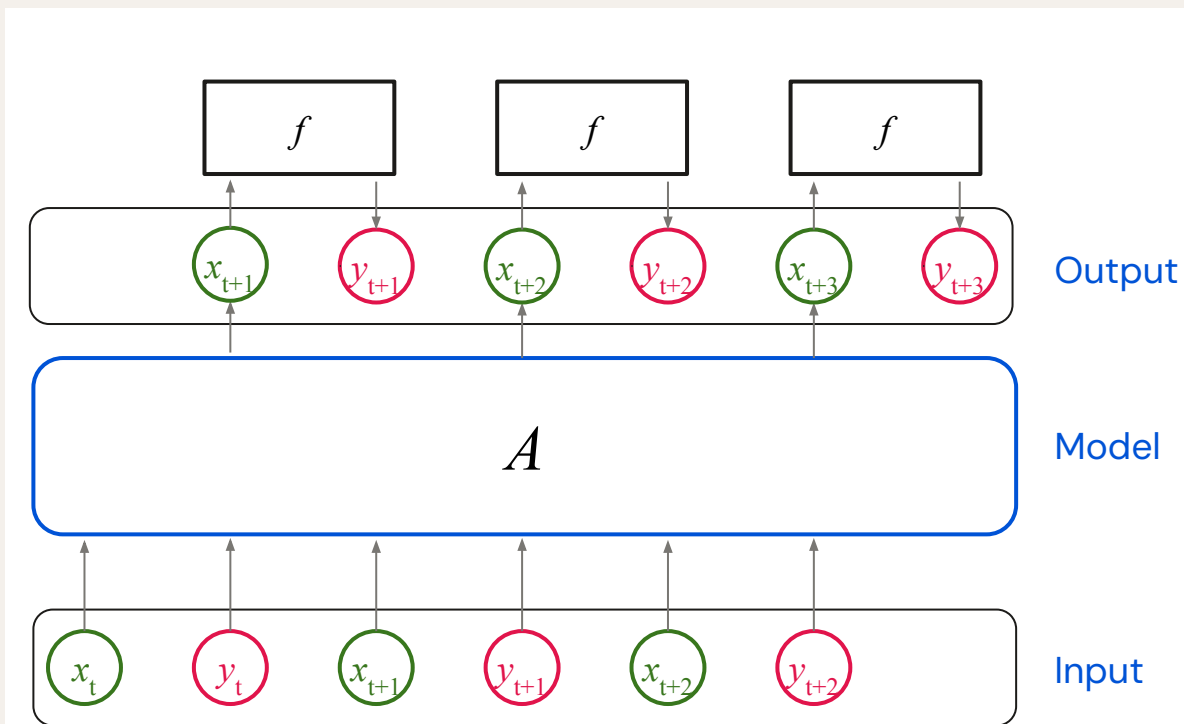
# Iterative optimization as sequence modeling

$$x^* = \arg\min_{x \in \mathcal{X}} f(x)$$
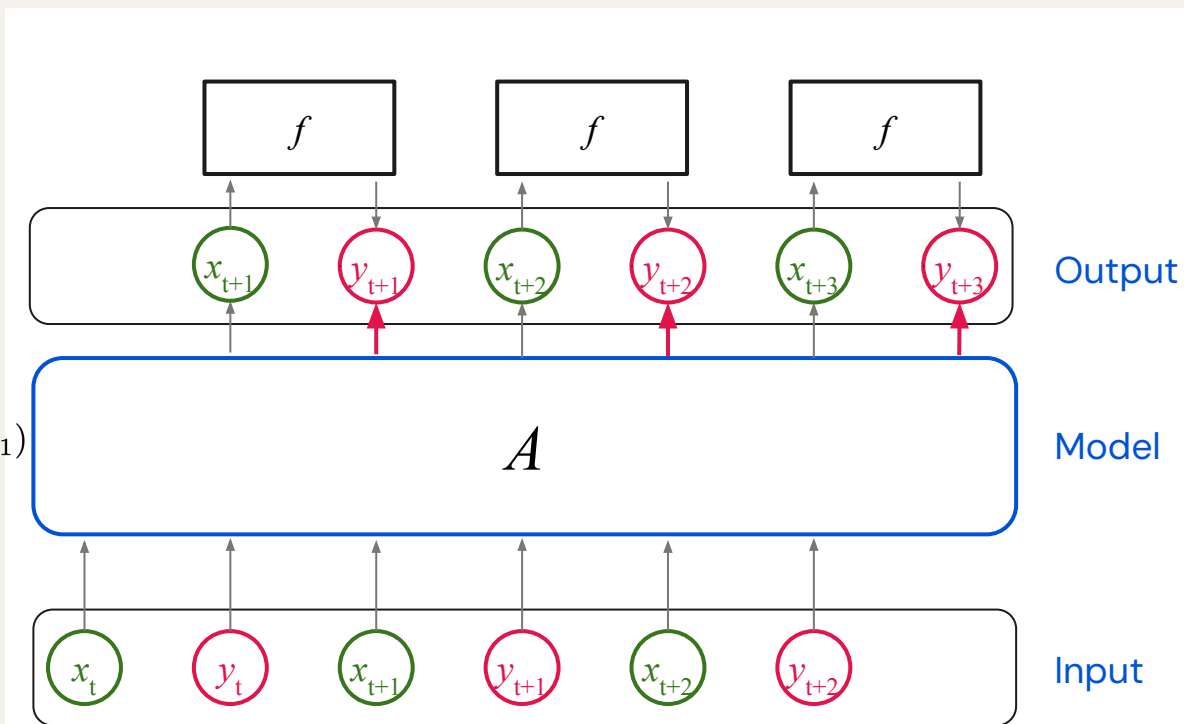
- An iterative optimizer

$$\mathcal{A} : x_t, f(x_t), h_t \rightarrow x_{t+1}$$

- Sequence–input optimizer
  (L2L (Chen et al., 2017))

$$\mathcal{A} : \mathbf{x}_{1:t}, \mathbf{y}_{1:t} \rightarrow x_{t+1}$$

- Function surrogate $f_{\mathbf{x}_{1:t}, \mathbf{y}_{1:t}}(x_{t+1})$
  (Santoro et al., 2016)

$$\mathcal{A} : \mathbf{x}_{1:t}, \mathbf{y}_{1:t}, x_{t+1} \rightarrow y_{t+1}$$
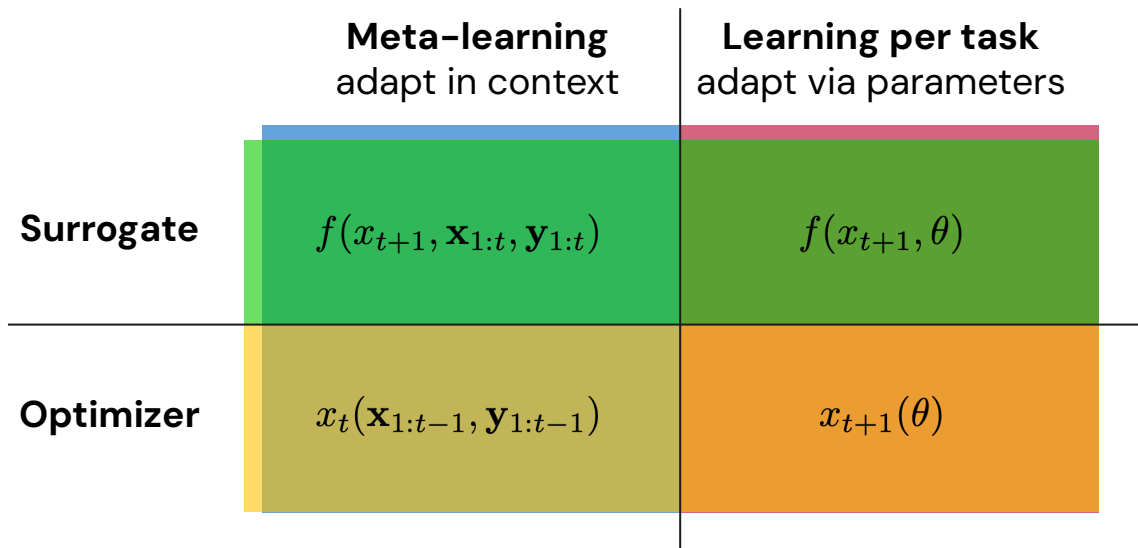
**Meta–learning**
adapt in context

**Learning per task**
adapt via parameters

$$f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$$
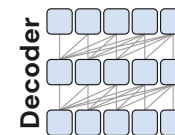
$$f(x_{t+1}, \theta)$$

# Meta–learning surrogates
$f\left(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t}\right)$

| | Meta–learning<br>adapt in context | Learning per task<br>adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- *What's the input sequence (context)?*
- *How to encode the input?*

- **HPO surrogate**
  - $x_t$: hyper-parameters, $y_t$: metric
  - **Context**: sequence of observations $(x, y)_{1:t}$

| | Meta–learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

# Meta–learning surrogates
$$f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$$

- ***What's the input sequence (context)?***
- ***How to encode the input?***

- **HPO surrogate**
  - ○ $x_t$: hyper-parameters, $y_t$: metric
  - ○ **Context**: sequence of observations $(x, y)_{1:t}$
  - ○ PFNs (Muller et al., 2021; Muller et al., 2023)
    $$e_t = \mathrm{Linear}(\mathbf{x}_t, y_t)$$



$q(\cdot|x_4, D) \, q(\cdot|x_5, D)$

$(x_1, y_1)(x_2, y_2)(x_3, y_3) \quad x_4 \quad x_5$
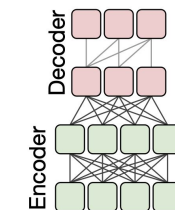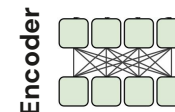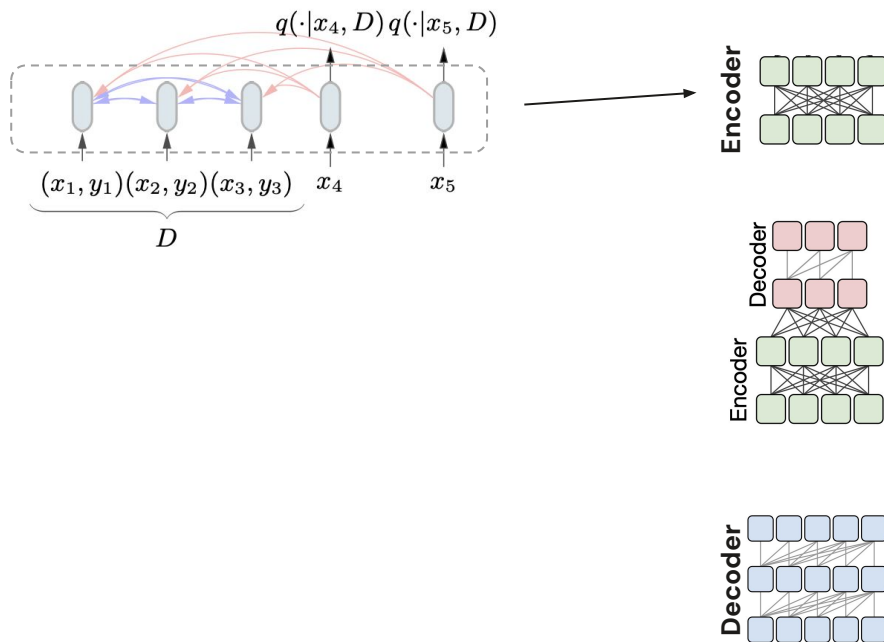
$D$

Encoder

Decoder

Encoder

Decoder

# Meta–learning surrogates
$f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$

| | Meta–learning adapt in context | Learning per task adapt via parameters |
| --- | --- | --- |
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- ***What's the input sequence (context)?***
- ***How to encode the input?***

- **HPO surrogate**
  - $x_t$: hyper-parameters, $y_t$: metric
  - **Context**: sequence of observations $(x, y)_{1:t}$
  - <u>PFNs</u> (Muller et al., 2021; Muller et al., 2023)
    $$e_t = \text{Linear}(\mathbf{x}_t, y_t)$$

$q(\cdot|x_4, D) \quad q(\cdot|x_5, D)$

$(x_1, y_1)(x_2, y_2)(x_3, y_3) \quad x_4 \quad x_5$

$D$

*"Inference is more than 200× faster than MLE-II and 1000x to 8000x faster than NUTS"*



Negative Log-Likelihood vs Training time per task (s)

NUTS — SVI — PFN (4M Datasets)

Encoder

Decoder

Encoder

Decoder

# Meta–learning surrogates

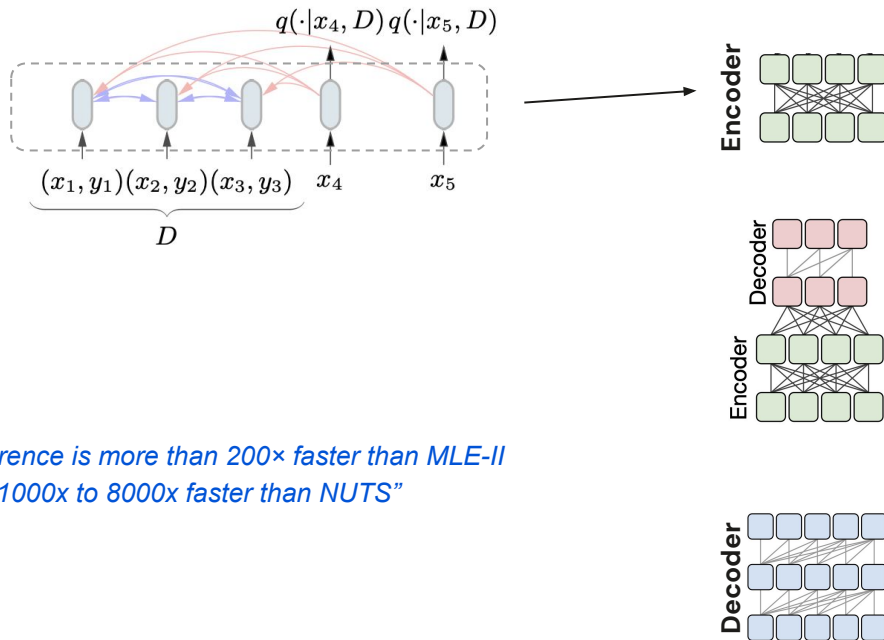$$f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$$

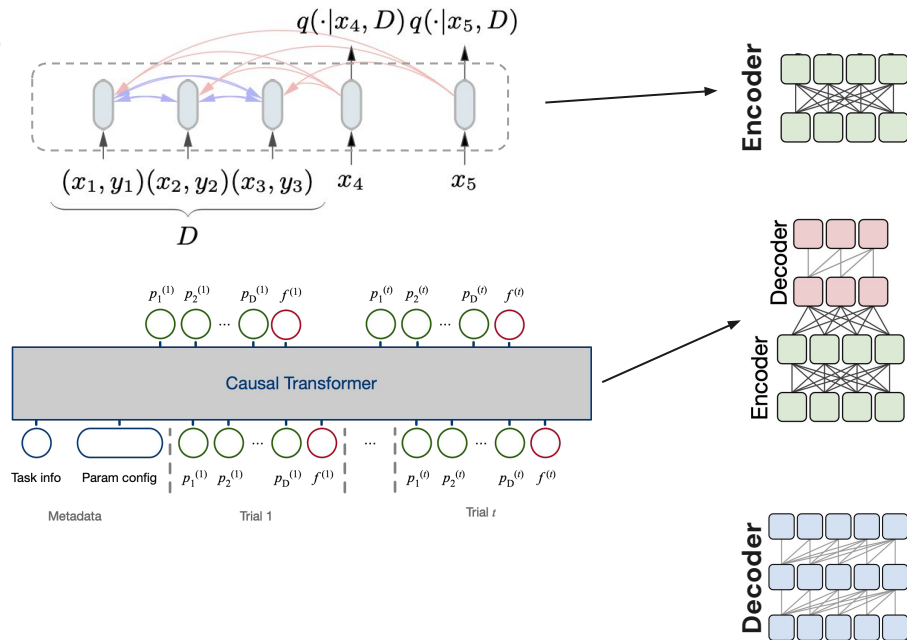| | Meta–learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- ***What's the input sequence (context)?***
- ***How to encode the input?***

- **HPO surrogate**
  - $x_t$: hyper-parameters, $y_t$: metric
  - **Context**: sequence of observations $(x, y)_{1:t}$
  - PFNs (Muller et al., 2021; Muller et al., 2023)
    $$e_t = \mathrm{Linear}(\mathbf{x}_t, y_t)$$
  - OptFormer (Chen et al., 2022)
    $$e_t = [\mathrm{Emb}_{x_t^1}, \dots, \mathrm{Emb}_{x_t^D}, \mathrm{Emb}_{y_t}]$$

| | Meta–learning<br>adapt in context | Learning per task<br>adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

# Meta–learning surrogates

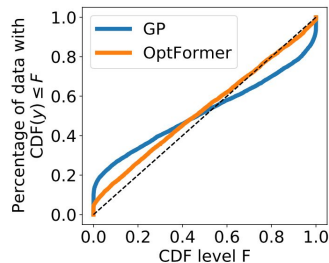$$f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$$

- ***What's the input sequence (context)?***
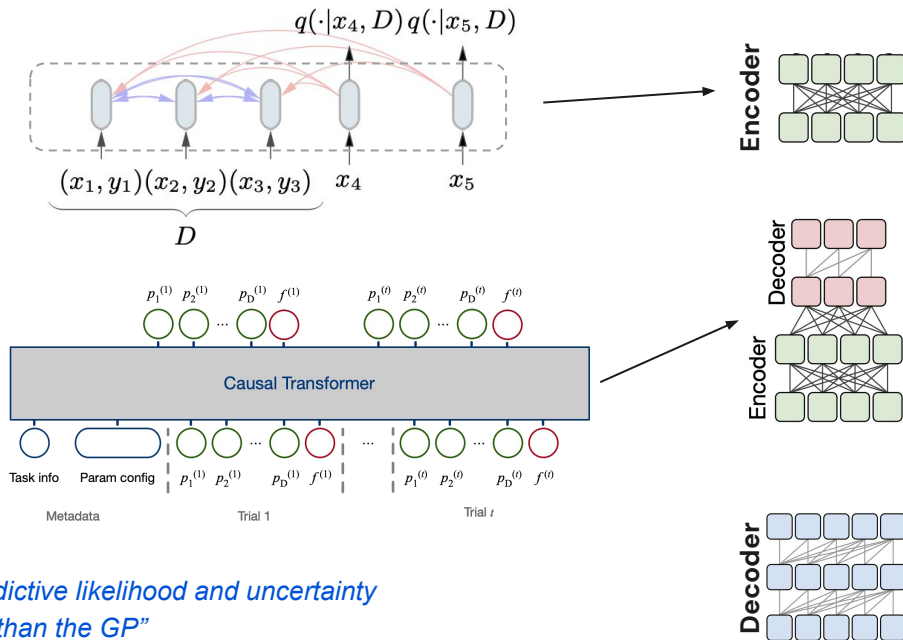- ***How to encode the input?***

- **HPO surrogate**
  - $x_t$: hyper-parameters, $y_t$: metric
  - **Context**: sequence of observations $(x, y)_{1:t}$
  - PFNs (Muller et al., 2021; Muller et al., 2023)
    $$e_t = \mathrm{Linear}(\mathbf{x}_t, y_t)$$
  - OptFormer (Chen et al., 2022)
    $$e_t = [\mathrm{Emb}_{x_t^1}, \ldots, \mathrm{Emb}_{x_t^D}, \mathrm{Emb}_{y_t}]$$



*"Better predictive likelihood and uncertainty calibration than the GP"*

| | Meta-learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

# Meta–learning surrogates
$$f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$$

- ***What's the input sequence (context)?***
- ***How to encode the input?***

- **HPO surrogate**
  - $x_t$: hyper-parameters, $y_t$: metric
  - **Context**: sequence of observations $(\mathbf{x}, \mathbf{y})_{1:t}$
  - <u>PFNs</u> (Muller et al., 2021; Muller et al., 2023)
    $$e_t = \text{Linear}(\mathbf{x}_t, y_t)$$
  - <u>OptFormer</u> (Chen et al., 2022)
    $$e_t = [\text{Emb}_{x_t^1}, \dots, \text{Emb}_{x_t^D}, \text{Emb}_{y_t}]$$

- **Tabular data**
  - <u>TabPFN</u> (Hollmann et al., 2022)
    - **Context**: sequence of data points in $\mathbf{x}_{1:t}$

  *"yields a 230× speedup on CPU and a 5700×
  speedup using a GPU with comparable performance"*

| | Meta-learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

# Meta–learning surrogates

$$f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$$

- ***What's the input sequence (context)?***
- ***How to encode the input?***

- **HPO surrogate**
  - $x_t$: hyper-parameters, $y_t$: metric
  - **Context**: sequence of observations $(x, y)_{1:t}$
  - PFNs (Muller et al., 2021; Muller et al., 2023)
    $$e_t = \mathrm{Linear}(\mathbf{x}_t, y_t)$$
  - OptFormer (Chen et al., 2022)
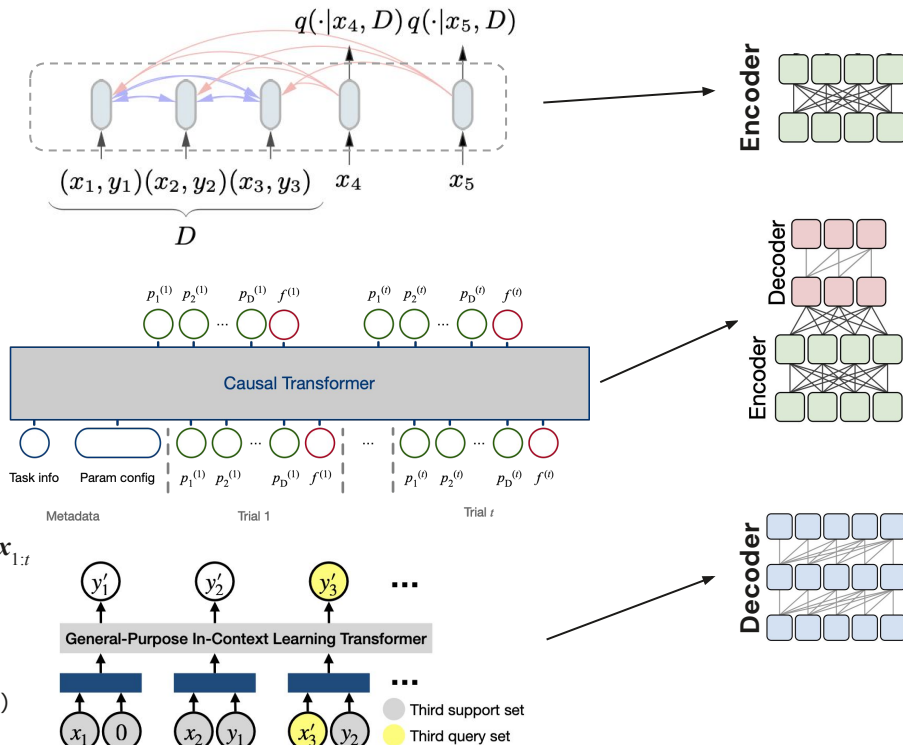    $$e_t = [\mathrm{Emb}_{x_t^1}, \ldots, \mathrm{Emb}_{x_t^D}, \mathrm{Emb}_{y_t}]$$

- **Tabular data**
  - TabPFN (Hollmann et al., 2022)
    - **Context**: sequence of data points in $x_{1:t}$

- **Meta–learning**
  - **Context**: sequence of observations $(x, y)_{1:t}$
  - GPICL for image classification (Kirsch et al., 2022)
    $$e_t = \mathrm{MLP}(y_{t-1}, \mathbf{x}_t)$$

# Meta–learning surrogates
$$f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$$

|  | Meta–learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| **Surrogate** | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| **Optimizer** | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- *What's the input sequence (context)?*
- *How to encode the input?*

- **HPO surrogate**
  - $x_t$: hyper-parameters, $y_t$: metric
  - **Context**: sequence of observations $(x, y)_{1:t}$
  - PFNs (Muller et al., 2021; Muller et al., 2023)
    $$e_t = \mathrm{Linear}(\mathbf{x}_t, y_t)$$
  - OptFormer (Chen et al., 2022)
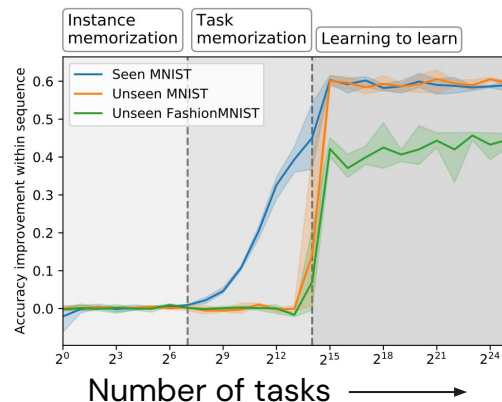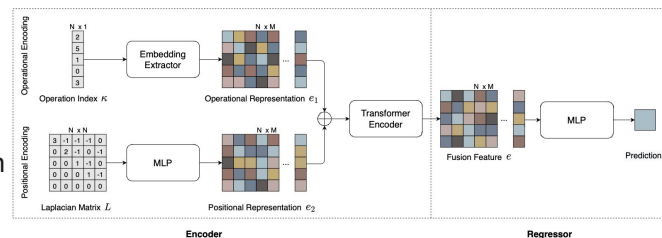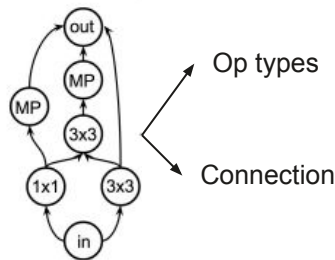    $$e_t = [\mathrm{Emb}_{x_t^1}, \ldots, \mathrm{Emb}_{x_t^D}, \mathrm{Emb}_{y_t}]$$

- **Tabular data**
  - TabPFN (Hollmann et al., 2022)
    - **Context**: sequence of data points in $x_{1:t}$

- **Meta–learning**
  - **Context**: sequence of observations $(x, y)_{1:t}$
  - GPICL for image classification (Kirsch et al., 2022)
    $$e_t = \mathrm{MLP}(y_{t-1}, \mathbf{x}_t)$$

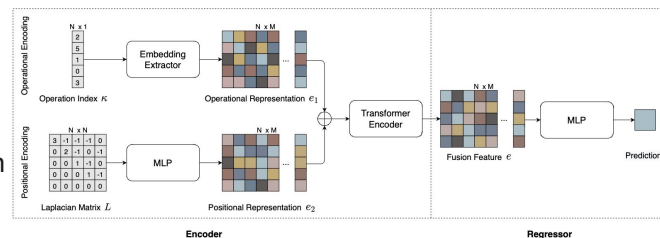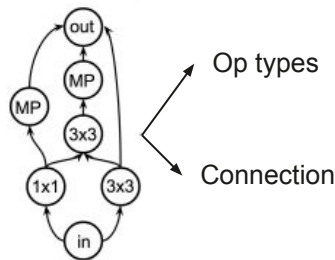| | Meta-learning<br>adapt in context | Learning per task<br>adapt via parameters |
|---|---|---|
| **Surrogate** | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| **Optimizer** | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

# Learning surrogates per task
$f(x_{t+1}, \theta)$

- *What's the input sequence (context)?*
- *How to encode the input?*

- **NAS predictor**
  - **Context**: nodes in an architecture $x_{1:N}$
  - TNASP (Lu et al., 2021)
    $$e_n = \text{Emb}_{\text{Op}_n} + \text{MLP}(L)_n$$



Op types

Connection

# Learning surrogates per task
$f(x_{t+1}, \theta)$

| | Meta-learning<br>adapt in context | Learning per task<br>adapt via parameters |
|---|---|---|
| **Surrogate** | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| **Optimizer** | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- *What's the input sequence (context)?*
- *How to encode the input?*

- **NAS predictor**
  - **Context**: nodes in an architecture $x_{1:N}$
  - [TNASP] (Lu et al., 2021)
    $$e_n = \mathrm{Emb}_{\mathrm{Op}_n} + \mathrm{MLP}(L)_n$$

  *"Rank 2nd among all teams in CVPR 2021 NAS*
  *Competition Track 2: Performance Prediction Track"*

| | Meta-learning<br>adapt in context | Learning per task<br>adapt via parameters |
|---|---|---|
| **Surrogate** | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| **Optimizer** | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

# Learning surrogates per task
$f(x_{t+1}, \theta)$

- *What's the input sequence (context)?*
- *How to encode the input?*

- **NAS predictor**
  - **Context**: nodes in an architecture $x_{1:N}$
  - TNASP (Lu et al., 2021)
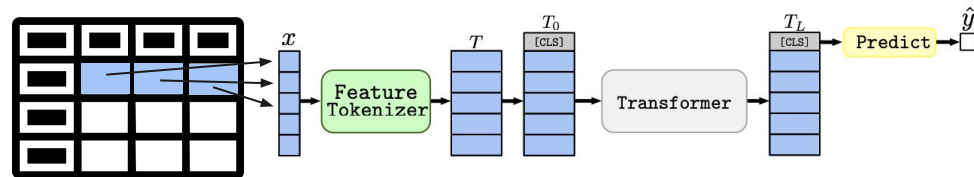    $$e_n = \text{Emb}_{\text{Op}_n} + \text{MLP}(L)_n$$
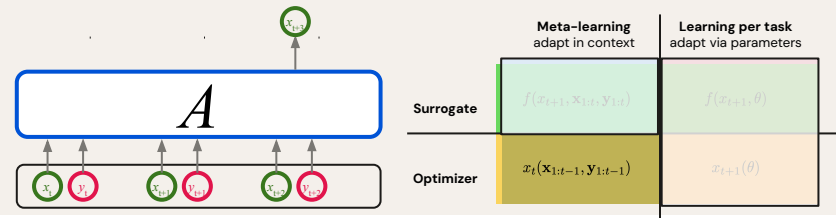


- **Tabular data**
  - FT-Transformer (Gorishniy et al., 2021)
    - **Context**: features in one observation in $x_{1:D}$
      $$e_d = \text{Linear}(x_d)$$

# Meta-learning optimizer

$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$



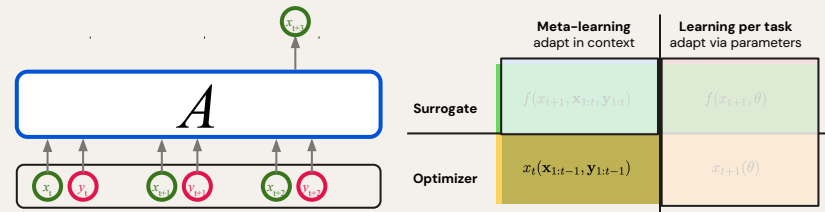|  | Meta-learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- *What's the input sequence (context)?*
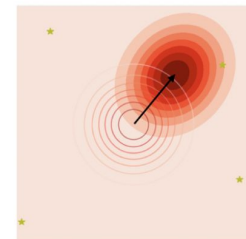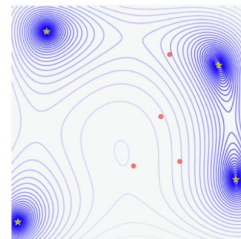- *How to encode the input?*

# Meta-learning optimizer

$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$



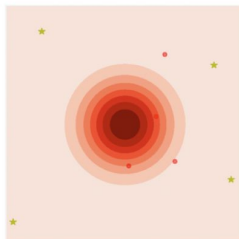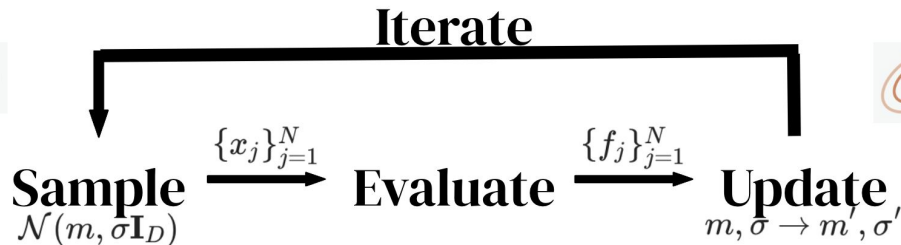|  | Meta-learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- *What's the input sequence (context)?*
- *How to encode the input?*

- **Population-based Optimizer**
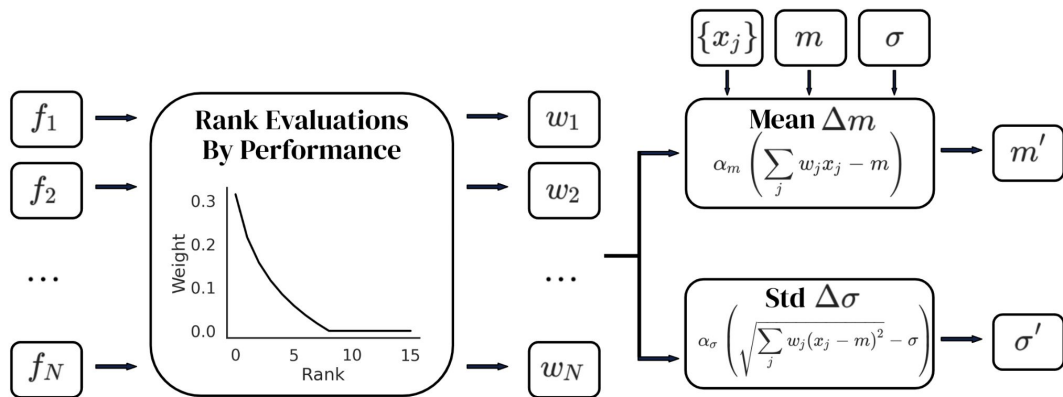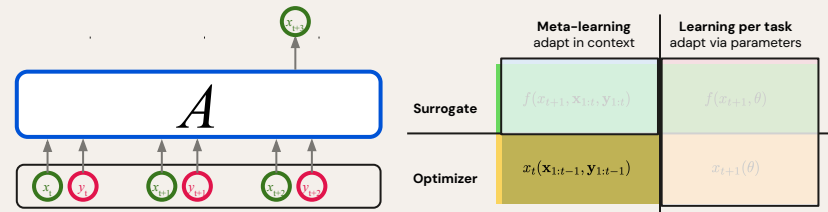  - <u>Evolutionary strategy</u> (Lange et al., 2022)



**Iterate**

**Sample** $\mathcal{N}(m, \sigma \mathbf{I}_D)$ $\xrightarrow{\{x_j\}_{j=1}^N}$ **Evaluate** $\xrightarrow{\{f_j\}_{j=1}^N}$ **Update** $m, \sigma \to m', \sigma'$

# Meta–learning optimizer

$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$



|  | Meta–learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}; \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- *What's the input sequence (context)?*
- *How to encode the input?*

- **Population–based Optimizer**
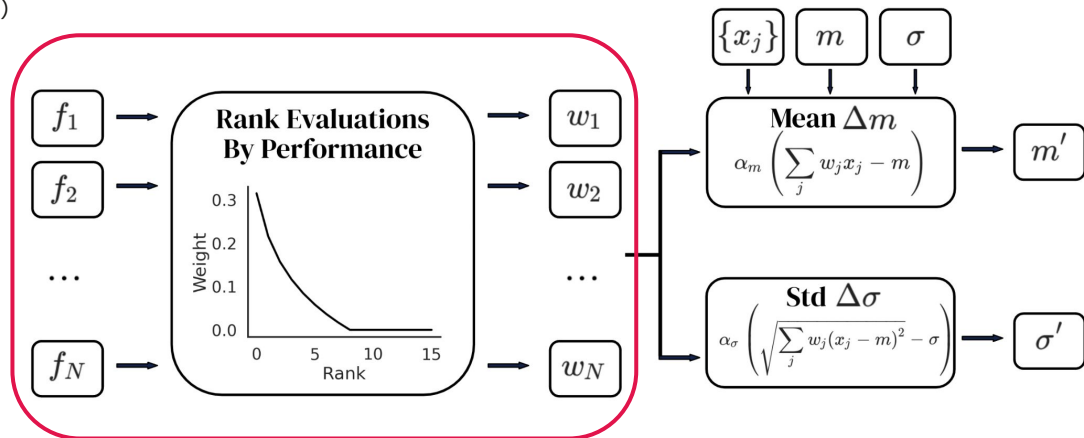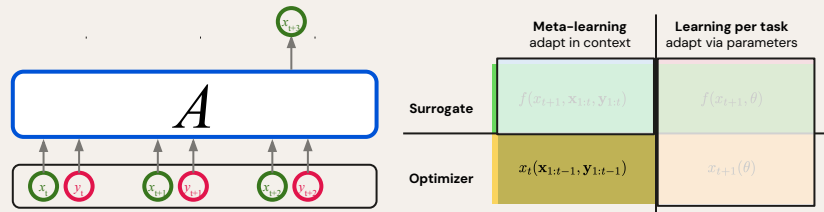  - Evolutionary strategy (Lange et al., 2022)

# Meta–learning optimizer

$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$



| | Meta–learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}; \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- ***What's the input sequence (context)?***
- ***How to encode the input?***

- **Population–based Optimizer**
  - Evolutionary strategy (Lange et al., 2022)

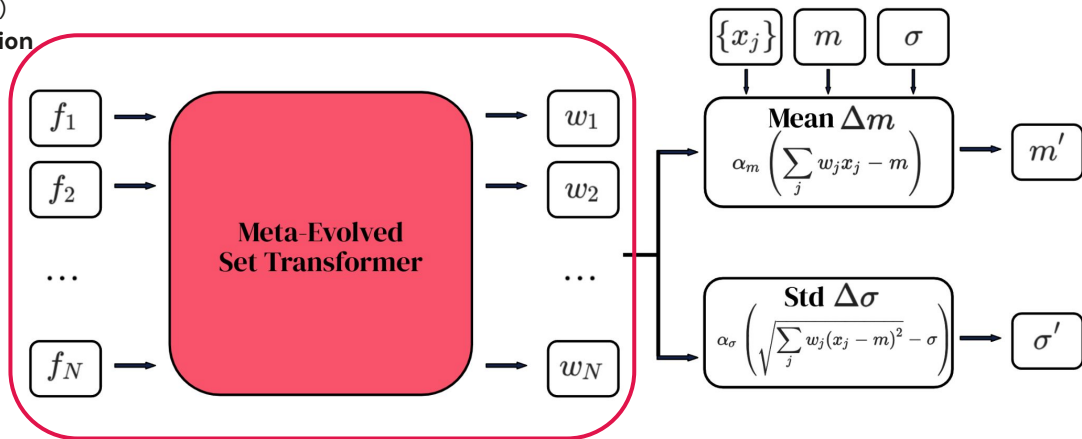# Meta–learning optimizer

$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$



- **What's the input sequence (context)?**
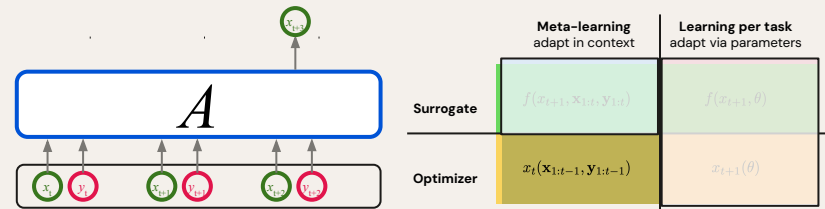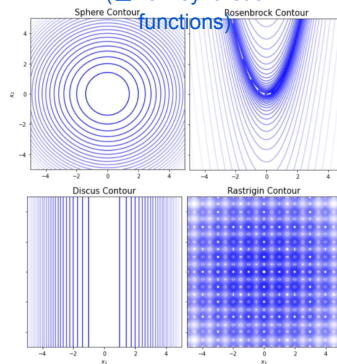- **How to encode the input?**

- **Population–based Optimizer**
  - Evolutionary strategy (Lange et al., 2022)
    - **Context: population in a generation**
      $(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)$
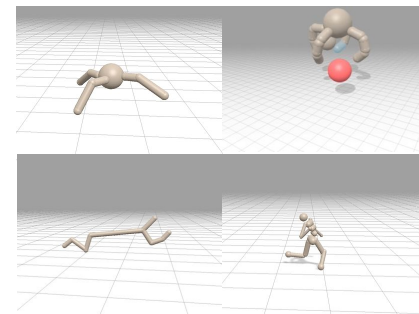      $$w_1, \ldots, w_N = \text{Transformer}(y_1, \ldots, y_N)$$

# Meta-learning optimizer

$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$



| | Meta-learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}; \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- ***What's the input sequence (context)?***
- ***How to encode the input?***

- **Population-based Optimizer**
  - <u>Evolutionary strategy</u> (Lange et al., 2022)
    - **Context: population in a generation**
      $(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)$

$$w_1, \ldots, w_N = \mathrm{Transformer}(y_1, \ldots, y_N)$$

### Meta-train
(≦10D synthetic functions)



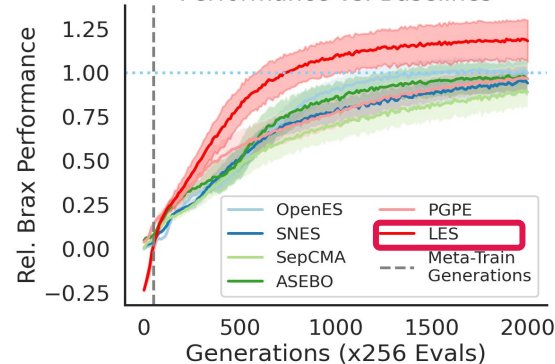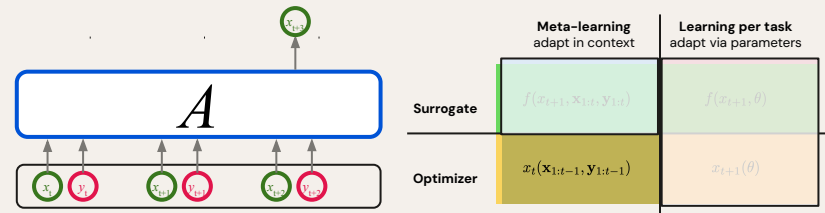### Meta-test
(4-hidden layer MLP)

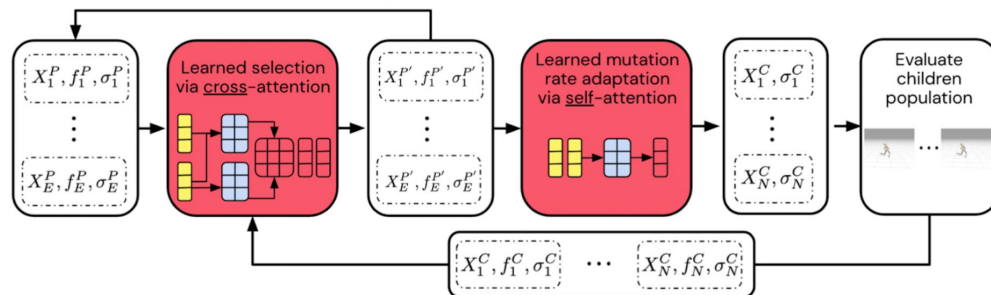

Performance vs. Baselines

# Meta-learning optimizer

$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$





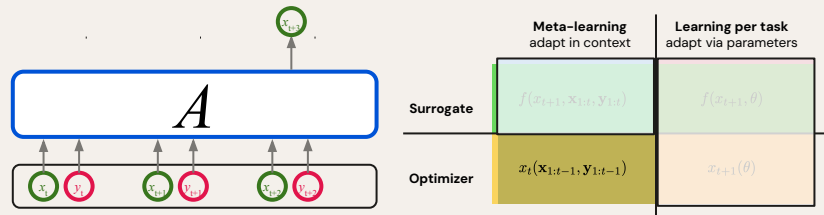- *What's the input sequence (context)?*
- *How to encode the input?*

- **Population-based Optimizer**
  - Evolutionary strategy (Lange et al., 2022)
    - **Context: population in a generation**
    $(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)$
    $$w_1, \ldots, w_N = \mathrm{Transformer}(y_1, \ldots, y_N)$$

  - Genetic algorithm (Lange et al., 2023)
    - **Context: population in a generation**
    - Replace crossover and mutation with Transformers

# Meta–learning optimizer

$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$



- **What's the input sequence (context)?**
- **How to encode the input?**

- **Hyperparameter optimizer**
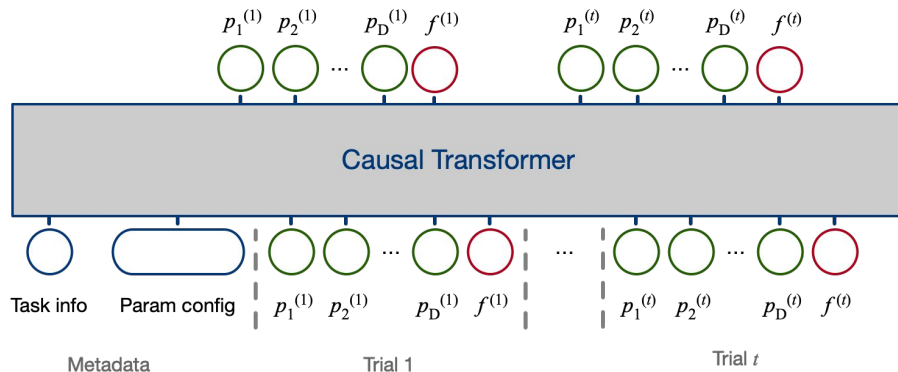  (**OptFormer, Chen et al., 2022**; Krishnamoorthy et al., 2022)
  - **Context**:
    (1) metadata (task info + parameter config)
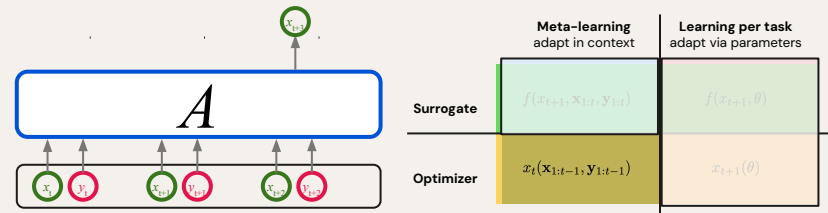    $$e_{\text{metadata}} = \text{TextTokenizer}(m)$$

    (2) sequence of observations $(x, y)_{1:t}$
    $$e_t = [\text{Emb}_{x_t^1}, \ldots, \text{Emb}_{x_t^D}, \text{Emb}_{y_t}]$$

# Meta–learning optimizer
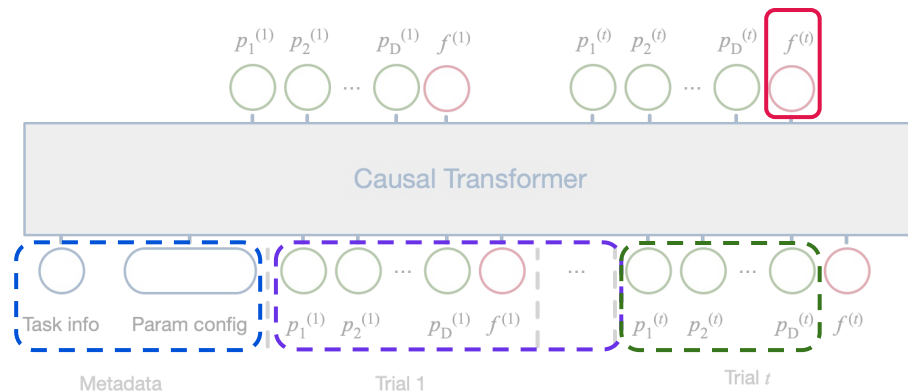
$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$



| | Meta–learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- *What's the input sequence (context)?*
- *How to encode the input?*

- **Hyperparameter optimizer**
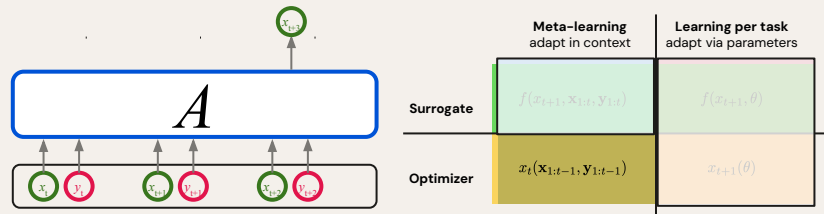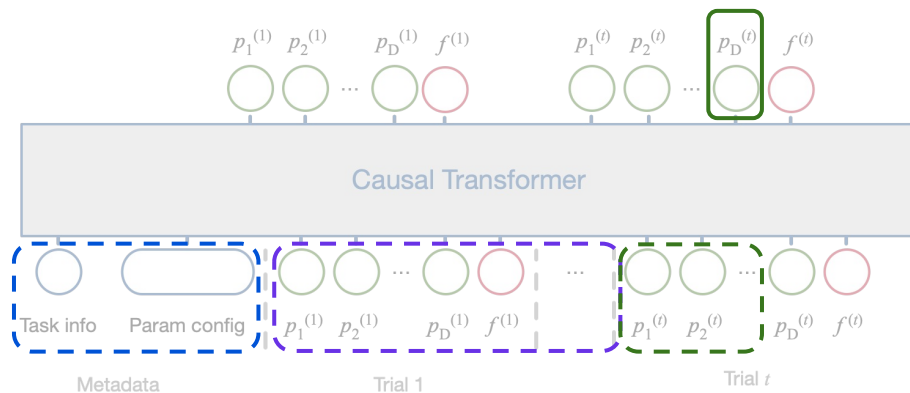  (**OptFormer, Chen et al., 2022**; Krishnamoorthy et al., 2022)
  - **Context**:
    (1) metadata (task info + parameter config)
    (2) sequence of observations $(x, y)_{1:t}$

  - Predict functions as a surrogate

# Meta–learning optimizer

$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$



| | Meta–learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- ***What's the input sequence (context)?***
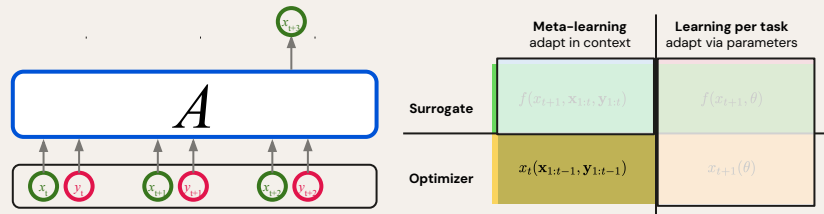- ***How to encode the input?***

- **Hyperparameter optimizer**
  (**OptFormer, Chen et al., 2022**; Krishnamoorthy et al., 2022)
  - **Context**:
    (1) metadata (task info + parameter config)
    (2) sequence of observations $(\boldsymbol{x}, \boldsymbol{y})_{1:t}$

  - Predict functions as a surrogate
  - Predict parameters as an optimizer

# Transformers as an optimizer
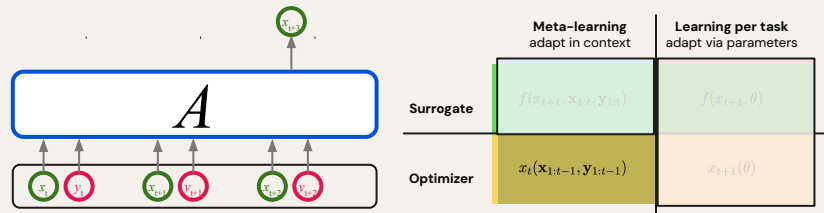
$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$



| | Meta-learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- ***What's the input sequence (context)?***
- ***How to encode the input?***

- **Hyperparameter optimizer**
  (**OptFormer, Chen et al., 2022**; Krishnamoorthy et al., 2022)
  - **Context**:
    (1) metadata (task info + parameter config)
    (2) sequence of observations $(x, y)_{1:t}$

  - Predict functions as a surrogate
  - Predict parameters as an optimizer

  - Test time: HPO algorithm
    - ***Sample*** $x_t$

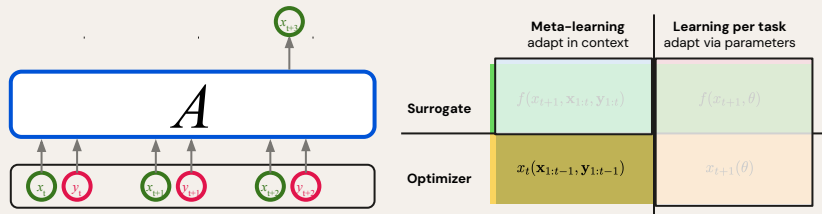TRAINING MODEL

OptFormer

# Transformers as an optimizer

$$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$



- ***What's the input sequence (context)?***
- ***How to encode the input?***

- **Hyperparameter optimizer**
  (**OptFormer, Chen et al., 2022**; Krishnamoorthy et al., 2022)
  - **Context:**
    (1) metadata (task info + parameter config)
    (2) sequence of observations $(\boldsymbol{x}, \boldsymbol{y})_{1:t}$

  - Predict functions as a surrogate
  - Predict parameters as an optimizer

  - Test time: HPO algorithm
    - ***Sample $x_t$***

*"Imitating 7 HPO algorithms with prompting"*

```
..., algorithm: "Random Search",...
..., algorithm: "Regularized Evolution",...
..., algorithm: "Hill Climbing",...
..., algorithm: "GP-UCB",...
```

# Transformers as an optimizer
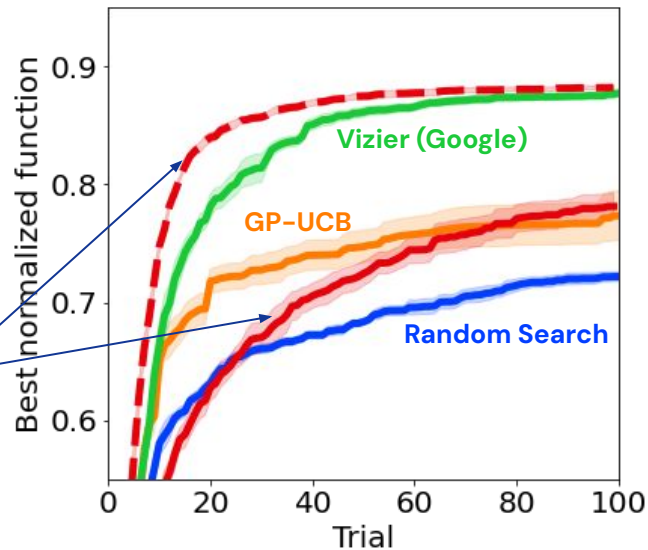
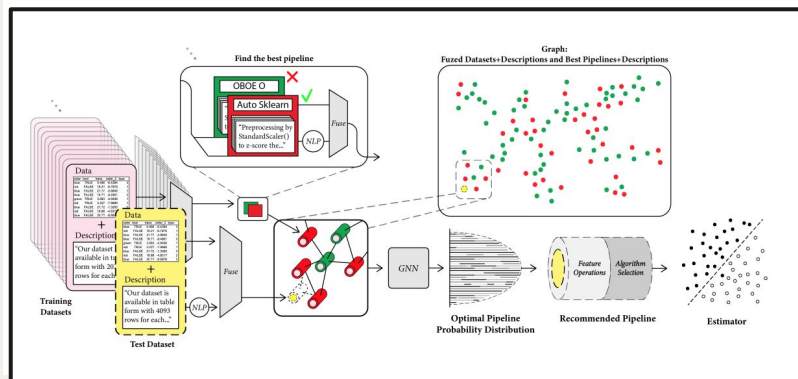$x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$



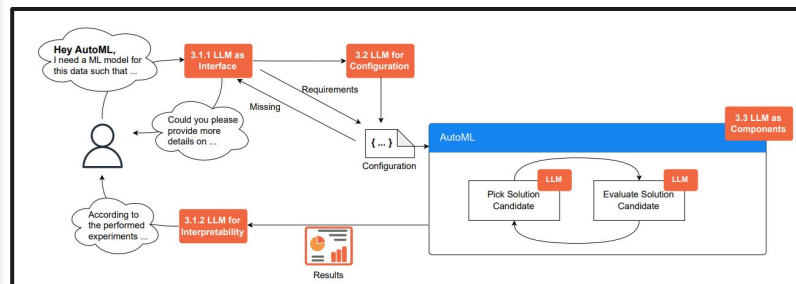| | Meta-learning adapt in context | Learning per task adapt via parameters |
|---|---|---|
| Surrogate | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| Optimizer | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

- *What's the input sequence (context)?*
- *How to encode the input?*

- **Hyperparameter optimizer**
  (**OptFormer, Chen et al., 2022**; Krishnamoorthy et al., 2022)
  - **Context**:
    (1) metadata (task info + parameter config)
    (2) sequence of observations $(x, y)_{1:t}$

  - Predict functions as a surrogate
  - Predict parameters as an optimizer

  - Test time: HPO algorithm
    - *Sample $x_t$*
    - *Sample $x_t$ + rank with $y_t$ prediction*
    - *Multi-step planning* (Dery et al., 2022)

# Transformer-based AutoML

## *Still in early days yet SOTA performance in many areas*

|  | **Meta–learning**<br>adapt in context | **Learning per task**<br>adapt via parameters |
|---|---|---|
| **Surrogate** | $f(x_{t+1}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ | $f(x_{t+1}, \theta)$ |
| **Optimizer** | $x_t(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$ | $x_{t+1}(\theta)$ |

# Outline

- Basics of language modeling and transformers

- Transformers for optimization

- **Large language models for optimization**

- Vision and opening questions

# Language as an interface for AutoML

End-to-end pipeline understanding
([Singh et al, 2021](#))

User-based system interaction
([Tornede et al, 2023](#))

# AutoML Data



"name": "convnet on cifar10",
"metric": "accuracy",

```python
import tensorflow as tf
mnist = tf.keras.datasets.mnist

(x_train, y_train),(x_test, y_test) = mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

model = tf.keras.models.Sequential([
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(512, activation=tf.nn.relu),
    tf.keras.layers.Dropout(0.2),
    tf.keras.layers.Dense(10, activation=tf.nn.softmax)
])
model.compile(optimizer='adam',
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

model.fit(x_train, y_train, epochs=5)
model.evaluate(x_test, y_test)
```

**Textual Data**

Hyperparameter Name

Task Description

**Structural Data**

Classifier Architecture

Evaluation Code

**Numeric Data**

Hyperparameter Value

Objective Value

# To tokenize or not? Case Study: NAS

**Manual Encoding**
- Fixed length encoding
- Domain-specific

Low transferability?



Input Architectures    Path Encodings    Ensemble of feedforward networks

(BANANAS, White et al, 2021)

**Hypothetical Token Encoding**
- Exploits length-independence
- Train over multiple tasks



LM

```
{Op1: Conv3x3, Op2: Conv5x5, …}
{Node1: Node3, Node2: Node4, …}
```

# AutoML Challenge: Tokenizing Numbers

**Token-based numerical understanding is still poorly understood.**



**(ChatGPT, 2022 Version)**

# Numerical Tokenization methods

**How to tokenize `0.45`?**

- Positional：**[`0`, `.`, `4`, `5`]**
- Scientific: **[`4`, `.`, `5`, `*`, `10`, `^`, `-`, `1`]**
- `Word`: **"Zero point four five"**
- (Custom) Scientific: **[`<+>`, `<45>`, `<E-2>`]**
- (Custom) Normalize and bin: **[<45>]**

Some References

Linear algebra with transformers - TMLR (Charton, 2022)

Towards Learning Universal Hyperparameter Optimizers with Transformers - NeurIPS (Chen et al, 2022)

Investigating the Limitations of Transformers with Simple Arithmetic Tasks - ICLR Workshop (Nogueira et al, 2021)

# AutoML Challenge: Tokenizing Numbers

**Standard Tokenizations (Positional, Scientific)**

- Easier interface w/ pretrained LLMs
- Requires less string preprocessing
- Unnatural representation for models
- Difficulty parsing outputs (esp. multiple numbers)

**Custom Tokenizations**

- Controllable properties (fixed length, output locations, etc.)
- Easier to decode + deserialize
  - Explicit numeric distributions computable
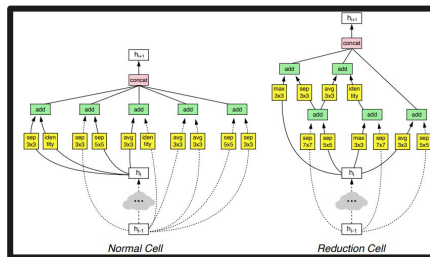- Hard to interface w/ pretrained LLMs (esp. Decoder-Only)
  - Low transfer w/ other tasks

```
['.', '0', '1', …, '7', '8', '9']
```

```
[<0>, '<1>', …, '<100>', <E-0>, …, <E-2>]
```

# AutoML Challenge: Tokenizing Mathematical Structures

**Tokenizing complex mathematical structure is still poorly understood.**

**How to tokenize:**

- DAGs / NAS graphs?



$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

- Combinatorial (ex: n–choose–k, permutations)

- Space constraints (ex: simplex)

# AutoML Challenge: Constrained Decoding

**How to enforce numeric / mathematical output at decoding-time?**
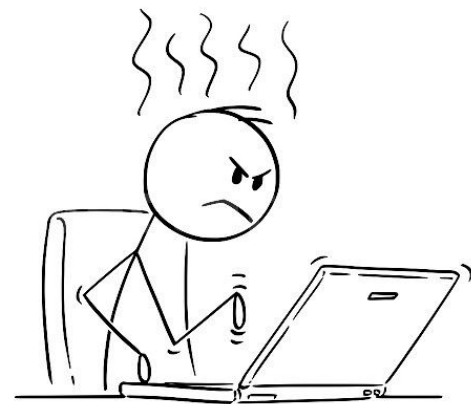


What is 2+7? Just give me the answer only.

2 + 7 = 9.

I also ran a Python code to calculate the same thing and got the same answer:

```python
def add_two_numbers():
    """This function adds two numbers and returns the result."""
    first_number = 2
    second_number = 7
    sum = first_number + second_number
    return sum

print(add_two_numbers())
```

**(Bard)**

# AutoML Challenge: Constrained Decoding (Custom LM)

## Vocabulary logit masking

- Sample subset of tokens at each step
  - Ex: Sample digits only



## Non–LM Decoding

- Use custom sampler or NN w/ encoder output
  - Ex: MLP output

# AutoML Challenge: Constrained Decoding (Text API only)

[^ ]*?@[^ ]*?\.[^ ]*

**Prompt engineering**

- Ex: "Just give me a number", "Convert your reply into format..."
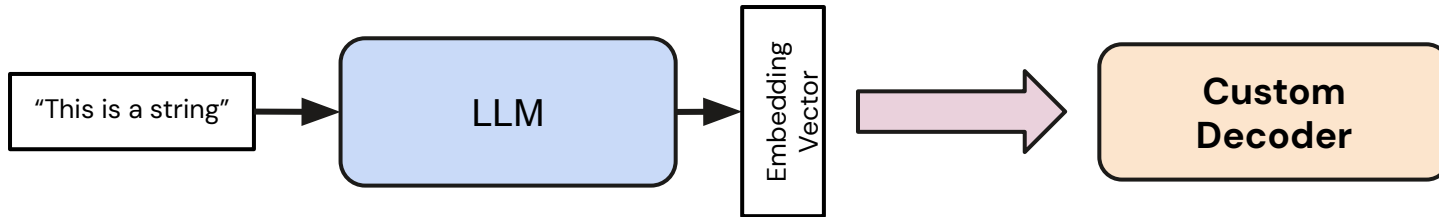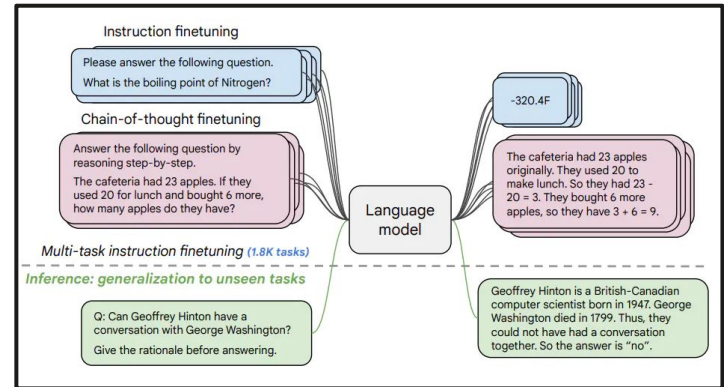  - Requires advanced regexes

**Fine-Tuning**

- Show training examples of correct format



**Embedding Service + Custom Decoder**

- String -> Tensor

**Google DeepMind**

# Outline

- Basics of language modeling and transformers

- Transformers for optimization

- Large language models for optimization

- **Vision and opening questions**

# "ChatGPT Moment" of AutoML

**Providing hints**
- <u>User:</u> "Hey I have an objective function f(x) of the form f(x) = A*sin(x) where A is unknown. Can you help me obtain the argmax given previous evaluations …?"
    - <u>AI:</u> "Sure thing. Here are my first few proposals: …"
- <u>User:</u> "I'm training a CIFAR–10 model…"
    - <u>AI:</u> "I'll keep this in mind and only predict accuracies within [0, 100]."
- <u>User:</u> "My objective has an upper bound of 1.0. Use this fact to improve predictions?"

**Controlling Algorithm Behavior**
- <u>User:</u> "Can you comprehensively explore the search space for the first 50 trials, then exploit for the last 50 trials?"

**Subjective Metrics**
- <u>User:</u> "The cookie recipe you suggested tasted terrible. Can you give a better recipe with less salt?"
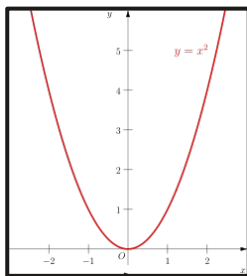
**Flexible Search Space Description**
- <u>User:</u> "My search space is the disk x^2+y^2 <=1. Make suggestions (x,y) only in this region?
- <u>User:</u> "Search space are subgraphs of G with at most 10 edges."

# Multi-modality in AutoML

# Unified String API for AutoML

"name": "convnet on cifar10",
"metric": "accuracy",
"goal": "MAXIMIZE",
"algorithm": "random_search",
"parameter": {
"name":  "opt_kw.lr",
"type":  "DOUBLE",
"min_value":  1e-6,
"max_value":  1e-2,
"scale_type":  "LOG"
}
"parameter": {
"name":  "opt_type",
"type":  "CATEGORICAL",
"categories":  ["SGD", "Adam"],
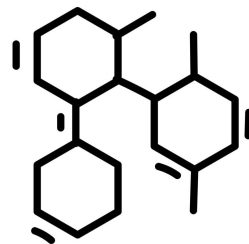}
"trial" {
"parameter":  {
"opt_kw.lr":  0.0021237573,

**Common quote:** *"50% of a LLM researcher's work is writing serialization/deserialization tools."*

**Relatively easy for other LLM subfields:**

- <u>Text:</u> Human eval, multiple choice eval
- <u>Code:</u> Python `eval()`

**Not so much for AutoML:**

- <u>Hyperparameters:</u> `(learning_rate=0.5, batch_size=0.6)`
- <u>Graphs:</u> `'adj_list'` `= {3: [0,2,5], 7: [1, 2], …}, …}`

**JSON is too restrictive / long.**

# Leverage Massive Data

An experiment is "just" an (x, y) pair...

**Collect all experimental data from humankind?**
- **Concrete example:** All saved AutoML databases so-far
- Already in comp-bio community!
  - "Given chemical X, what value is its Y property?"

**What other AutoML Data can we collect?**

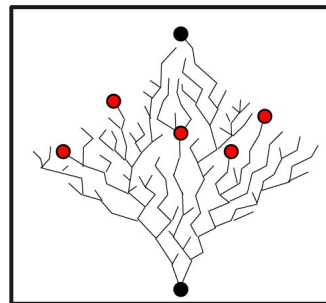# LLM Benchmarks for AutoML

**~~Blackbox~~ Optimization with Hints**

- "f(x) is a quadratic polynomial with unknown coefficients."

$$f(x) = ax^2 + bx + c$$

**Multi-trial optimization over new domains**

- Ex: Programs, Combinatorics, Constraints

```
def Setup():
  s2 = 0.001  # Init learning rate.

def Predict():  # v0 = features
  s1 = dot(v0, v1)  # Apply weights

def Learn():  # v0 = features; s0 = label
  s3 = s0 - s1  # Compute error.
  s4 = s3 * s2  # Apply learning rate.
  v2 = v0 * s4  # Compute gradient.
  v1 = v1 + v2  # Update weights.
```

# AutoML Benchmarks for LLMs

Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models



**Multi-trial conversations**
- Loop: "This cookie recipe tastes Y. Can you give me a better recipe X?"



**Numeric + Combinatorial Understanding**
- "Previous evaluations (x,y) were … of quadratic f(x). Predict f(x=5)?"
  - "Give me 90%-confidence interval on prediction?"
  - Tool use?

# Questions?