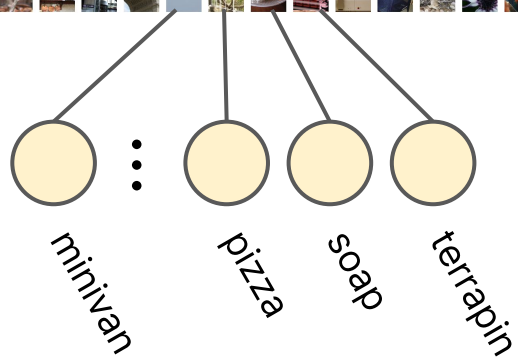
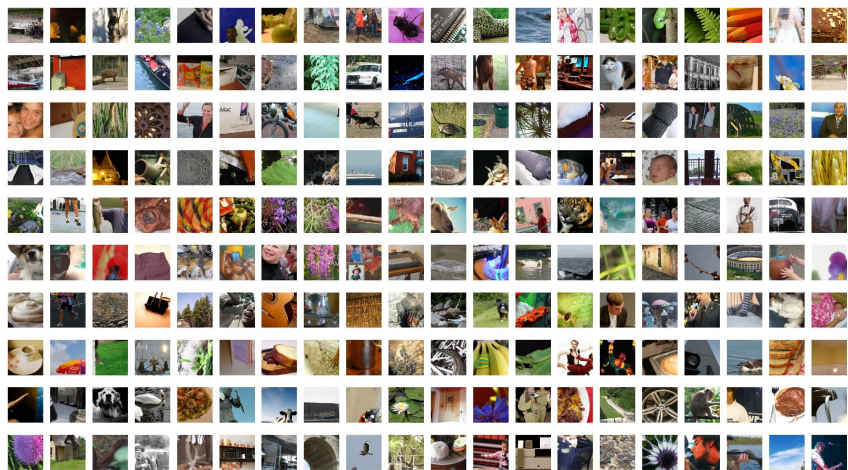


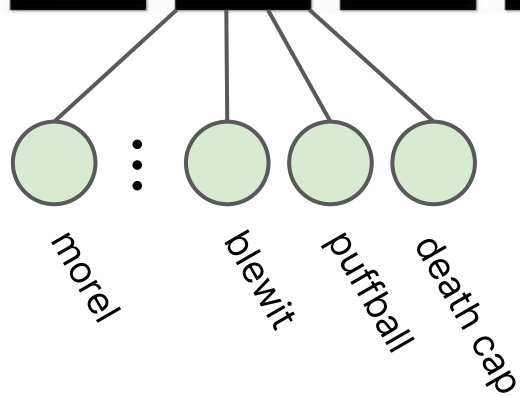
# A call to build models like we build open-source software

Colin Raffel

## Unsupervised pre-training



## Supervised fine-tuning



## *Unsupervised pre-training*

The cabs \_\_\_ the same rates as those \_\_\_ by horse-drawn cabs and were \_\_\_ quite popular; \_\_\_ the Prince of Wales (the \_\_\_ King Edward VII) travelled in \_\_\_. The cabs quickly \_\_\_ known as "hummingbirds" for \_\_\_ noise made by their motors and their distinctive black and \_\_\_ livery.

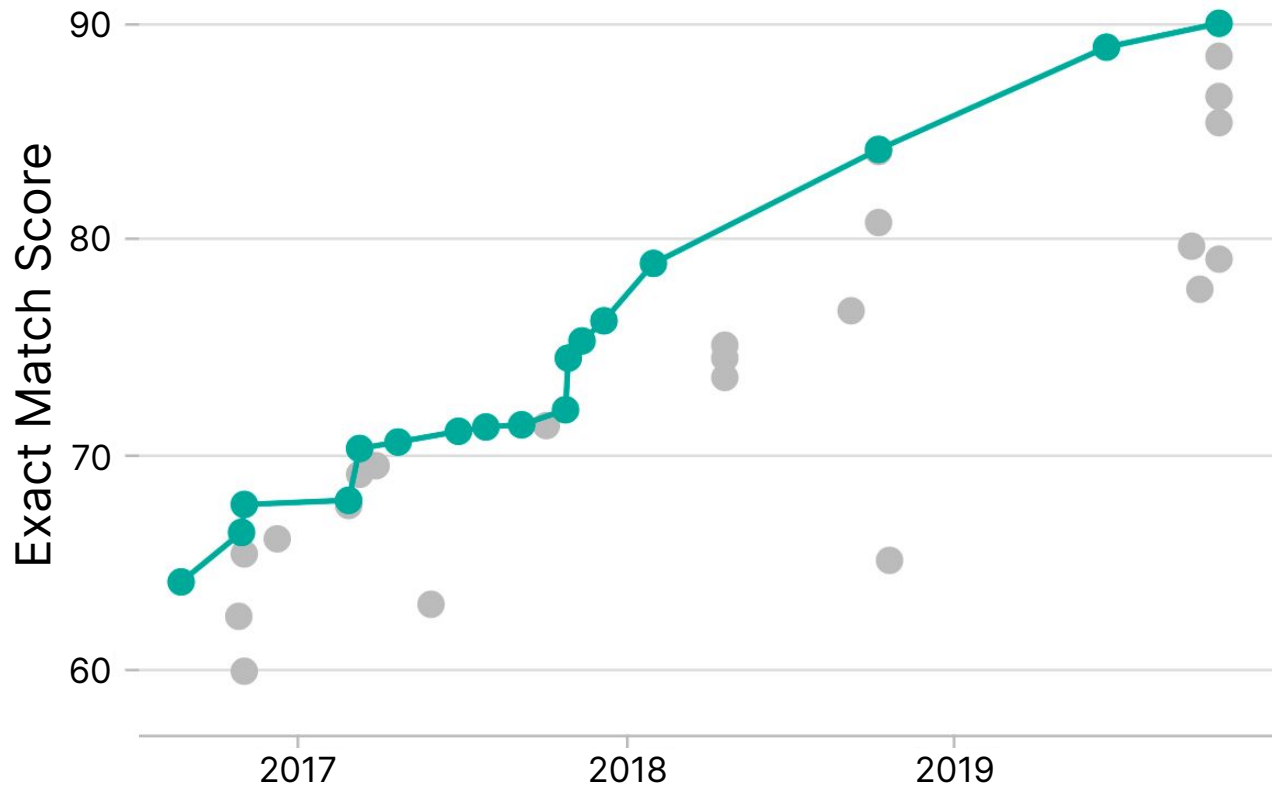
The cabs **charged** the same rates as those **used** by horse-drawn cabs and were **initially** quite popular; **even** the Prince of Wales (the **future** King Edward VII) travelled in **one**. The cabs quickly **became** known as "hummingbirds" for **the** noise made by their motors and their distinctive black and **yellow** livery.

## *Supervised fine-tuning*

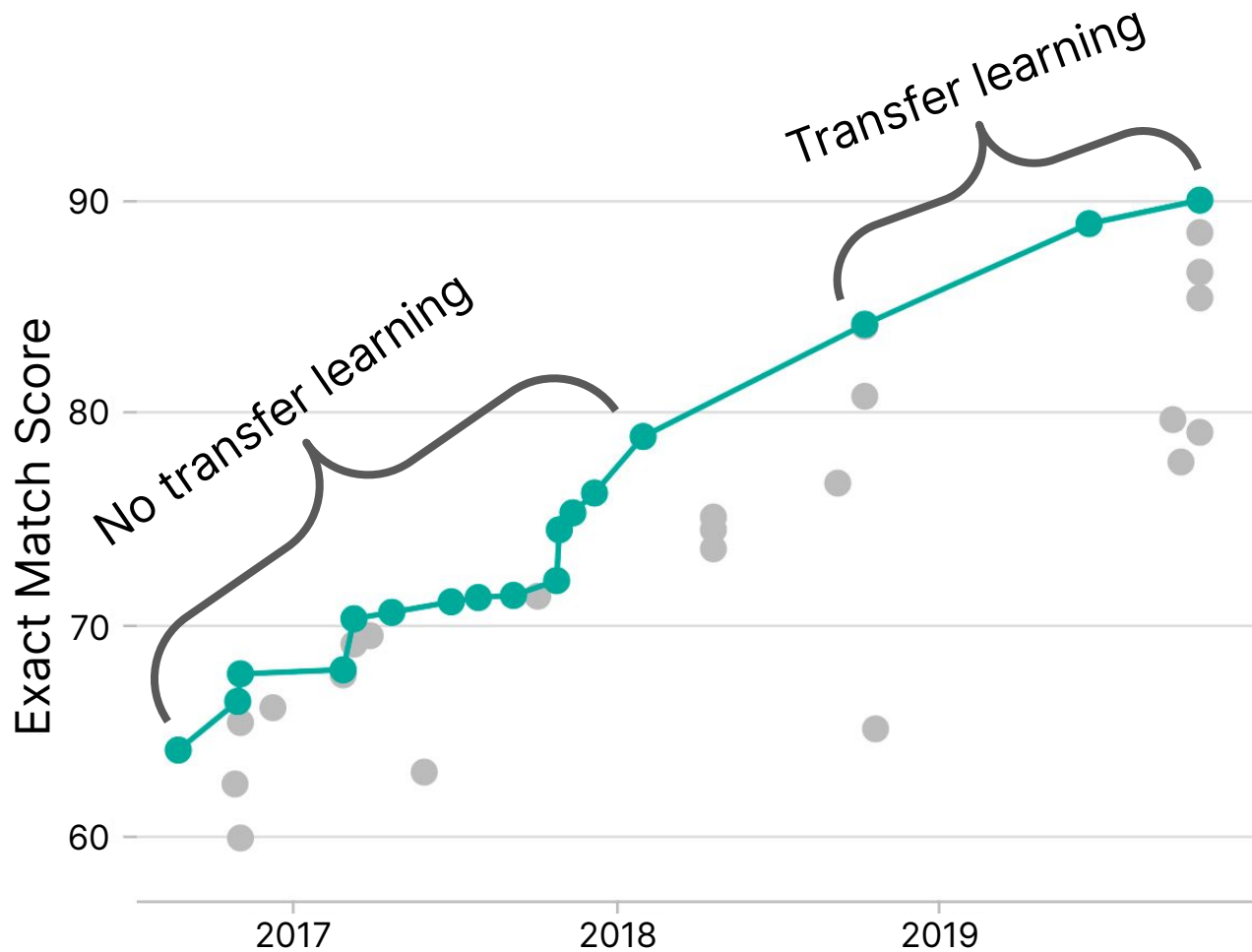
This movie is terrible! The acting is bad and I was bored the entire time. There was no plot and nothing interesting happened. I was really surprised since I had very high expectations. I want 103 minutes of my life back!

negative

# SQuAD Exact Match score

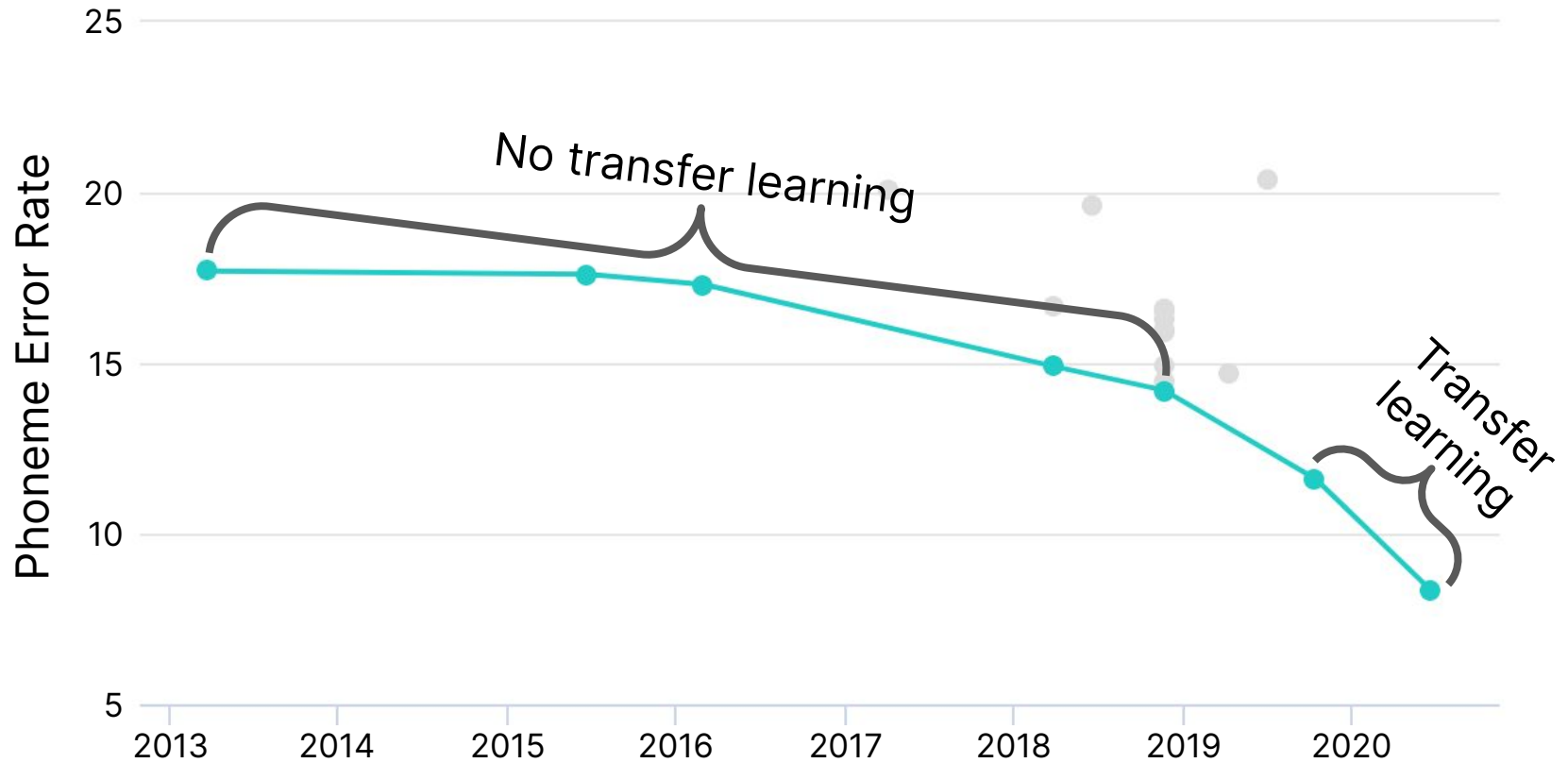


from <https://paperswithcode.com/sota/question-answering-on-squad11-dev>



from <https://paperswithcode.com/sota/question-answering-on-squad11-dev>

# TIMIT Phoneme Error Rate



from <https://paperswithcode.com/sota/speech-recognition-on-timit>



## TORCHVISION.MODELS

We provide pre-trained models, using the PyTorch [torch.utils.model\\_zoo](#). These can be constructed by passing `pretrained=True`:

```
import torchvision.models as models
resnet18 = models.resnet18(pretrained=True)
alexnet = models.alexnet(pretrained=True)
squeezenet = models.squeezenet1_0(pretrained=True)
vgg16 = models.vgg16(pretrained=True)
densenet = models.densenet161(pretrained=True)
inception = models.inception_v3(pretrained=True)
googlenet = models.googlenet(pretrained=True)
shufflenet = models.shufflenet_v2_x1_0(pretrained=True)
```

*GPT-3 175B model required  $3.14E23$  FLOPS of computing for training. Even at theoretical 28 TFLOPS for V100 and lowest 3 year reserved cloud pricing we could find, this will take 355 GPU-years and cost **\$4.6M** for a single training run.*





Models 33,490

Search Models

Add filters

Sort: Most Downloads



bert-base-uncased

Fill-Mask • Updated May 18 • ↓ 30M • ♥ 54



roberta-large

Fill-Mask • Updated May 21 • ↓ 13.1M • ♥ 20



distilbert-base-uncased

Fill-Mask • Updated Aug 29 • ↓ 4.83M • ♥ 26



xlm-roberta-base

Fill-Mask • Updated Sep 16 • ↓ 4.78M • ♥ 11



bert-base-cased

Fill-Mask • Updated Sep 6 • ↓ 4.02M • ♥ 6



distilbert-base-uncased-finetuned-sst-2-english

Text Classification • Updated Feb 9 • ↓ 3.54M • ♥ 18



roberta-base

Fill-Mask • Updated Jul 6 • ↓ 3.45M • ♥ 6



gpt2

Text Generation • Updated May 19 • ↓ 3.34M • ♥ 24

co:here

API

AI21 studio

# Custom language models built for scale

Build sophisticated language applications on top of AI21's language models

## Microsoft Megatron-Turing NLG 530B

The World's Largest and Most Powerful Generative Language Model

SambaNova<sup>®</sup>  
SYSTEMS

PRODUCTS ▾ SOLUTIONS ▾ RESOUR

Enterprise-Grade  
Large Language Models  
Made Simple & Accessible

Introducing Dataflow-as-a-Service™ GPT

OpenAI API Beta ABOUT EXAMPLES DOCS PRICING LOG IN [JOIN >](#)

## OpenAI technology, just an HTTPS call away

Apply our API to any language task — semantic search, summarization, sentiment analysis, content generation, translation, and more — with only a few examples or by specifying your task in English.

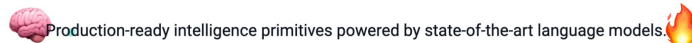
[JOIN THE WAITLIST >](#)

[</> EXPLORE THE DOCS](#)

Introducing the LightOn Muse API



# Create. Process. Understand. Learn.



Production-ready intelligence primitives powered by state-of-the-art language models.

For the first time natively in French, Spanish, Italian, and more. **Now in private beta!**

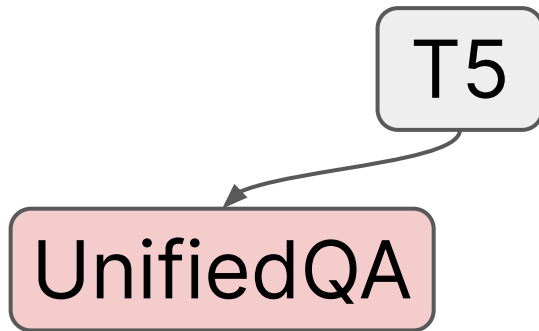
# AI, 모두의 능력이 된다. HyperCLOVA

AI가 모두의 능력이 되는 새로운 시대.

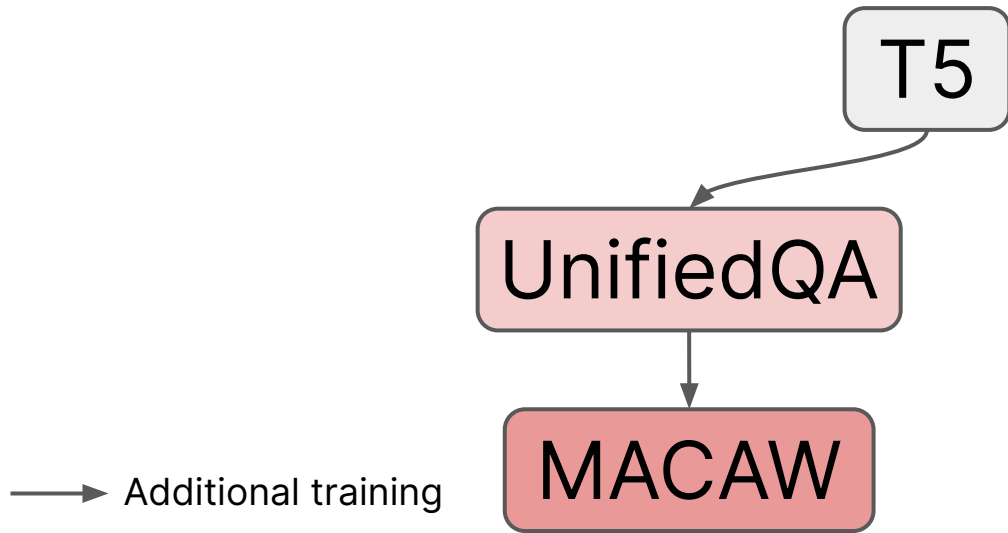
그 시작이 될 HyperCLOVA를 소개합니다.

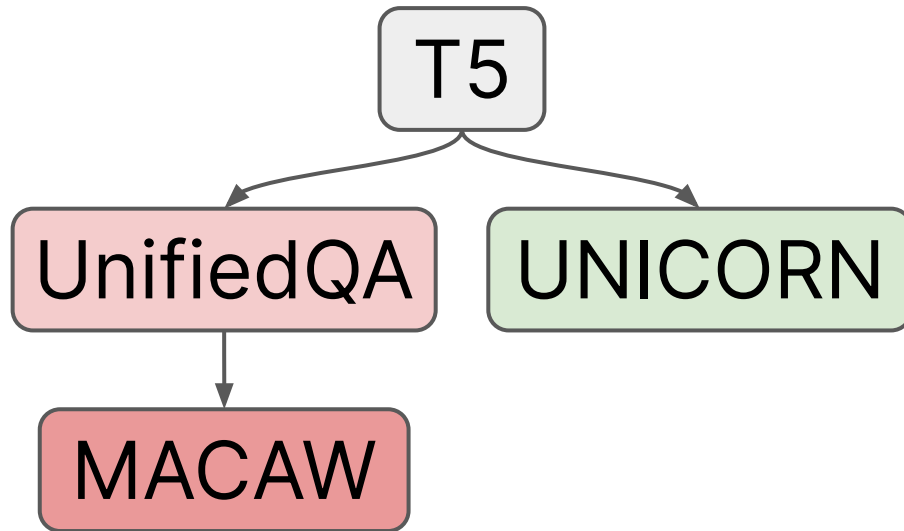
네이버 클로바와 함께 새로운 시대를 시작하세요.

T5

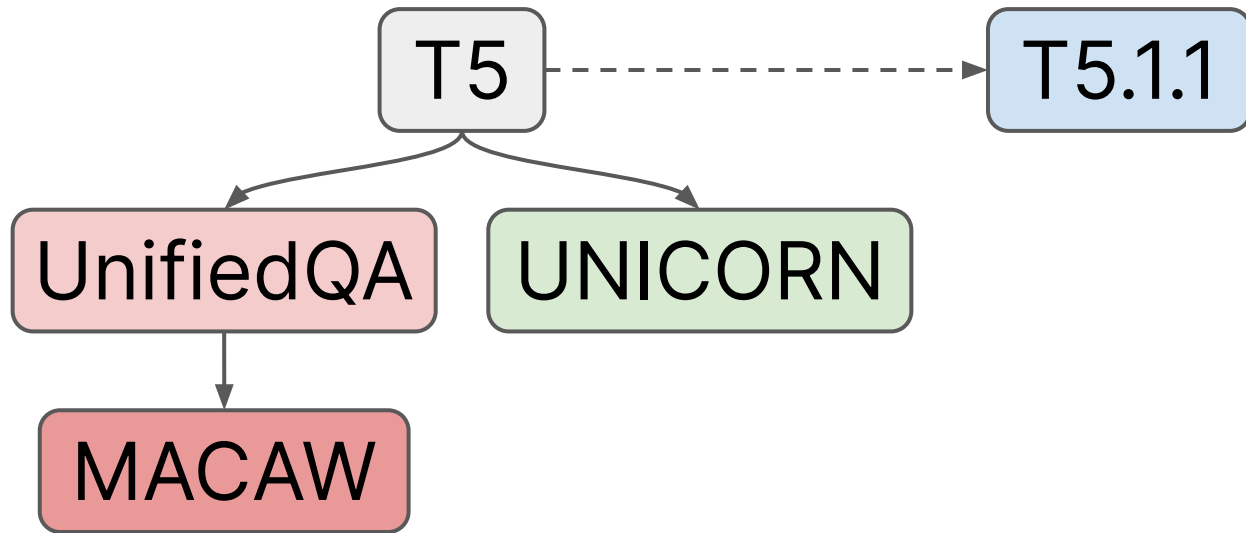


→ Additional training



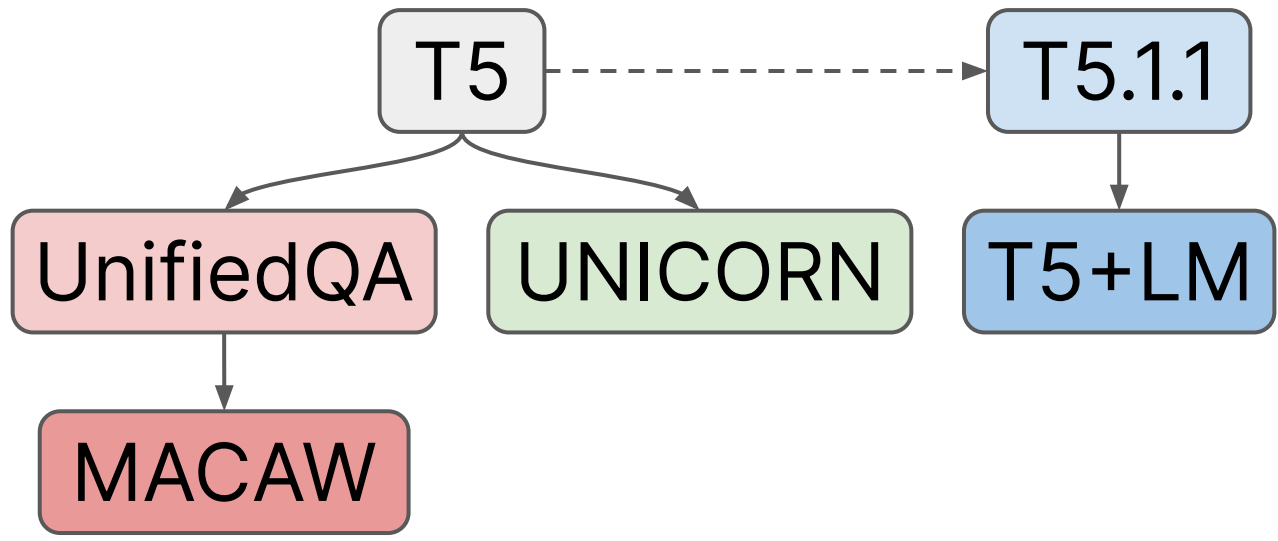


→ Additional training



—> Additional training

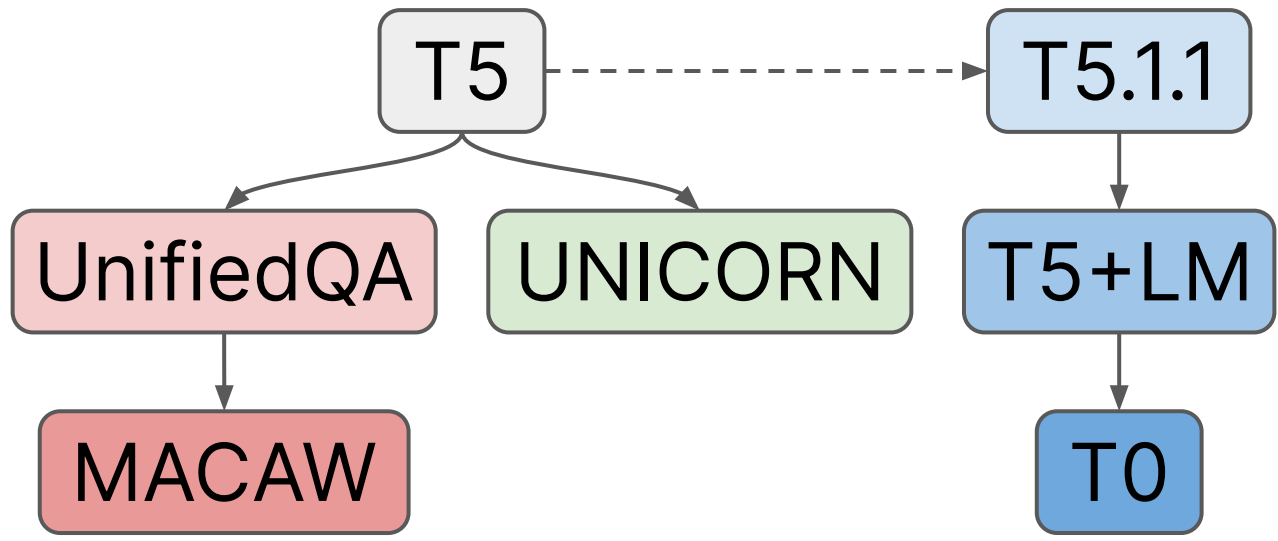
- - -> New model



—> Additional training

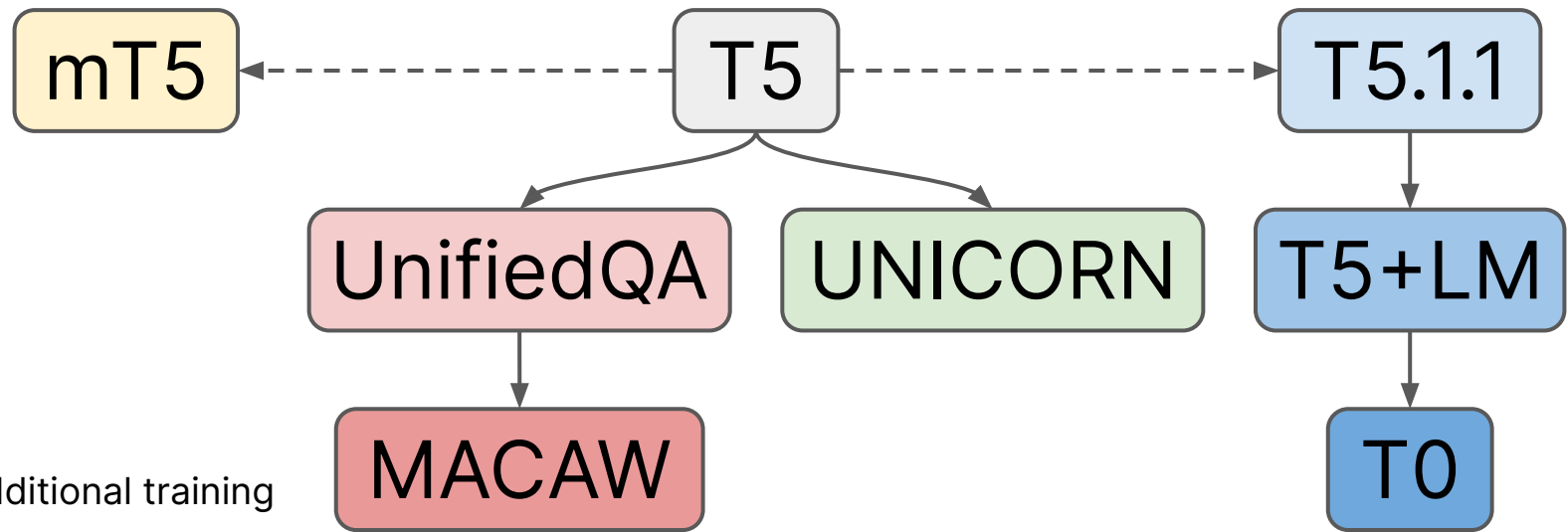
- - -> New model





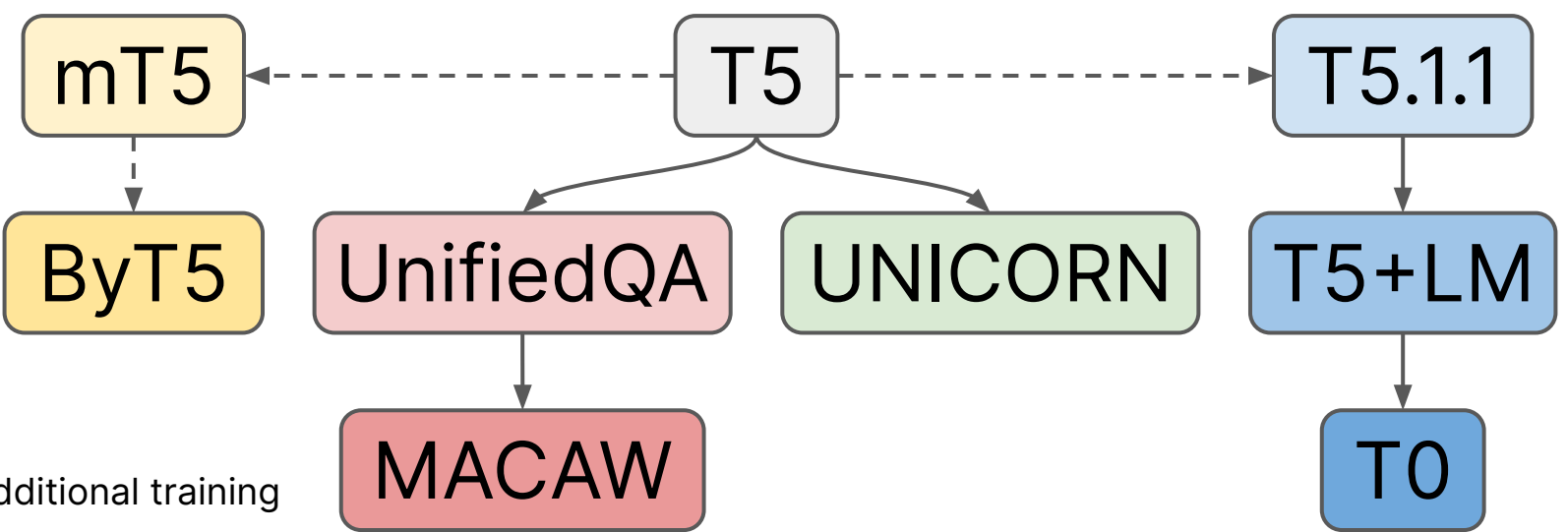
→ Additional training

- - - → New model



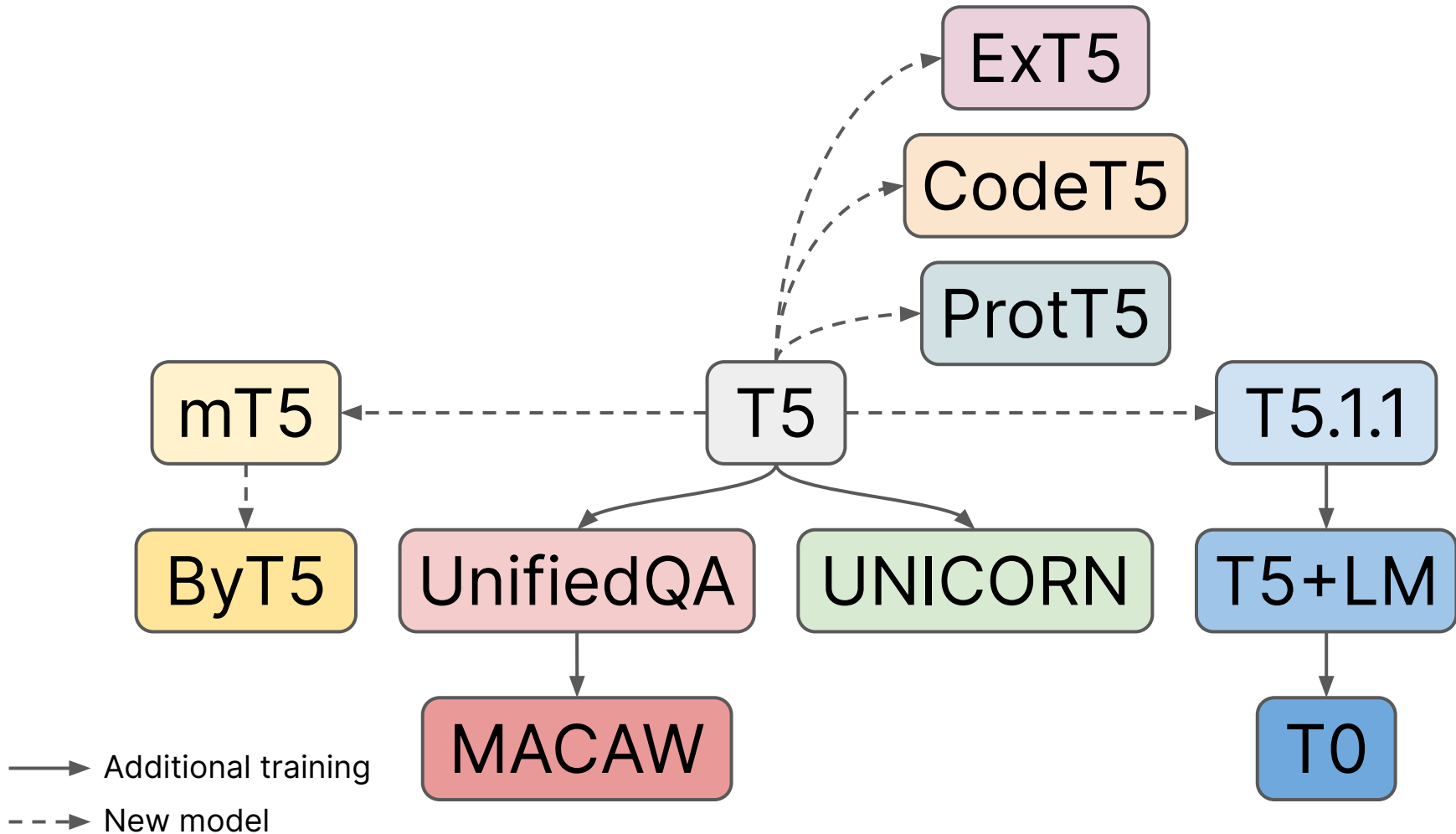
→ Additional training

- - - → New model



—> Additional training

- - -> New model



Models 1,407

t5

Add filters

Sort: Most Downloads

Vamsi/T5\_Paraphrase\_Paws  
Text Generation · Updated Jun 23 · ↓ 93.1k · ♥ 2

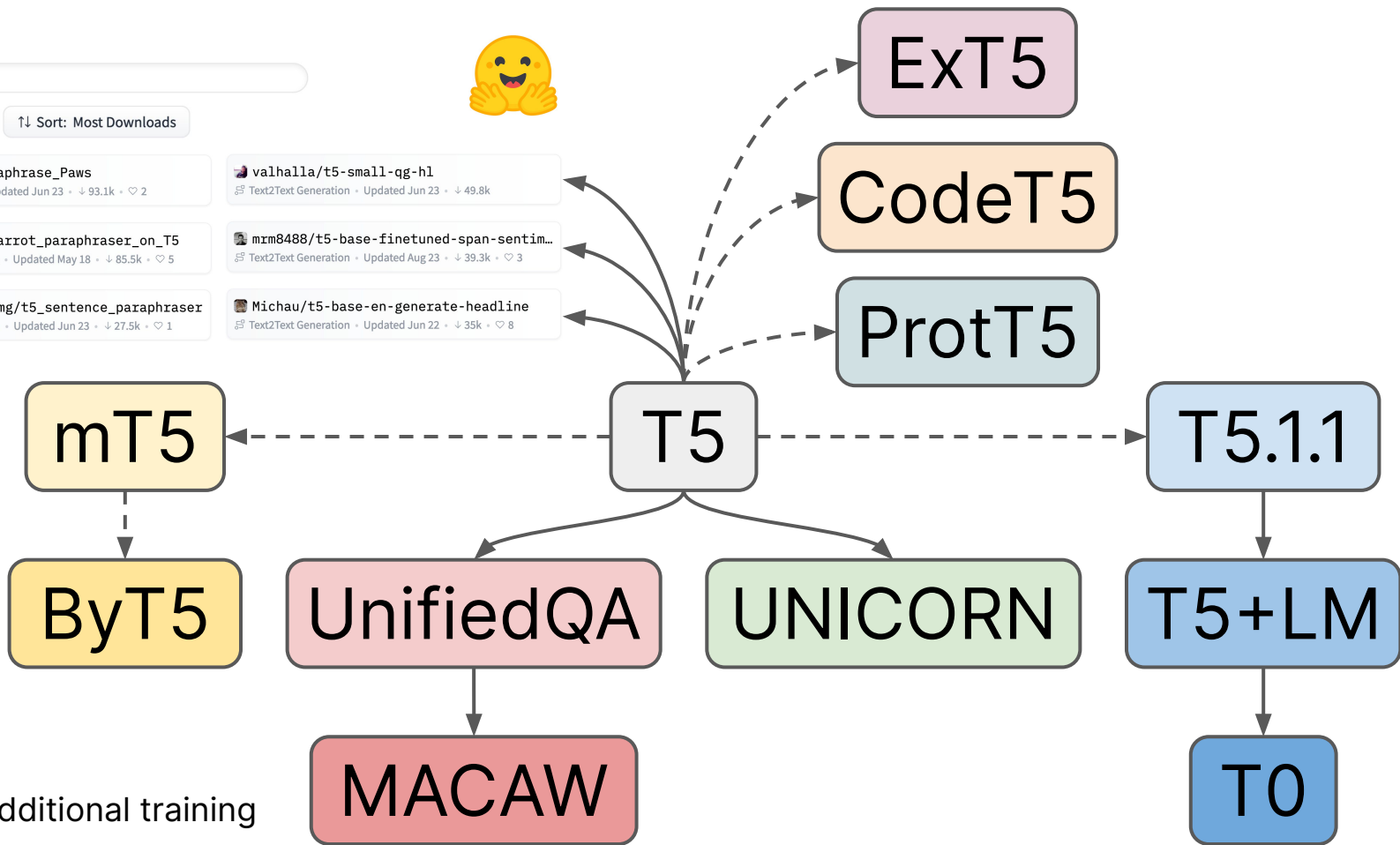
valhalla/t5-small-qg-h1  
Text2Text Generation · Updated Jun 23 · ↓ 49.8k

prithivida/parrot\_paraphraser\_on\_t5  
Text2Text Generation · Updated May 18 · ↓ 85.5k · ♥ 5

mrm8488/t5-base-finetuned-span-sentim...  
Text2Text Generation · Updated Aug 23 · ↓ 39.3k · ♥ 3

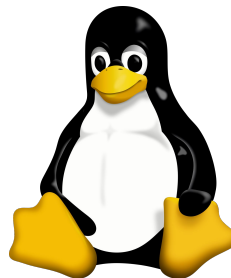
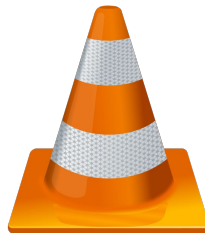
ramsrigoutham/t5\_sentence\_paraphraser  
Text2Text Generation · Updated Jun 23 · ↓ 27.5k · ♥ 1

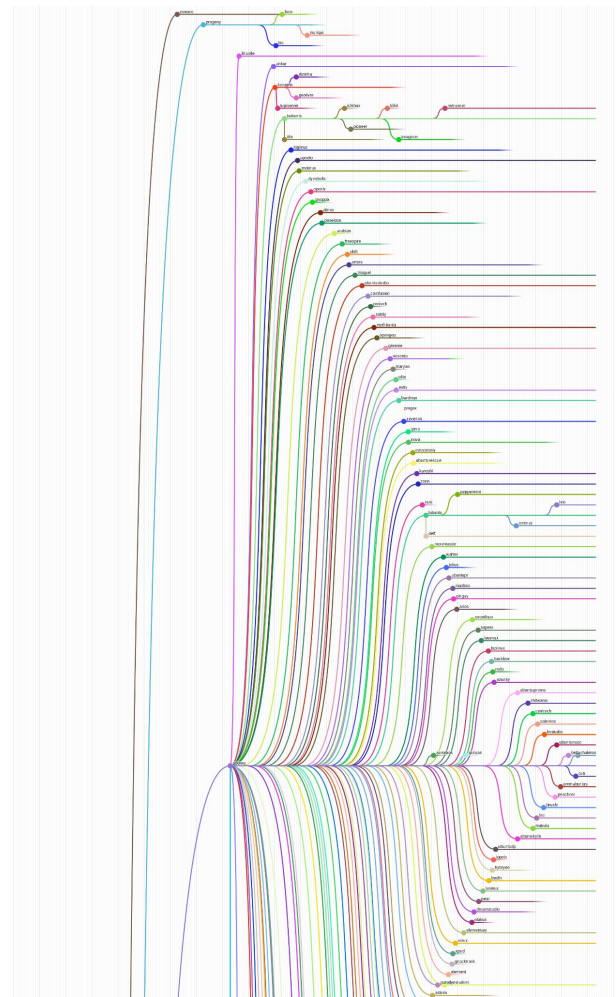
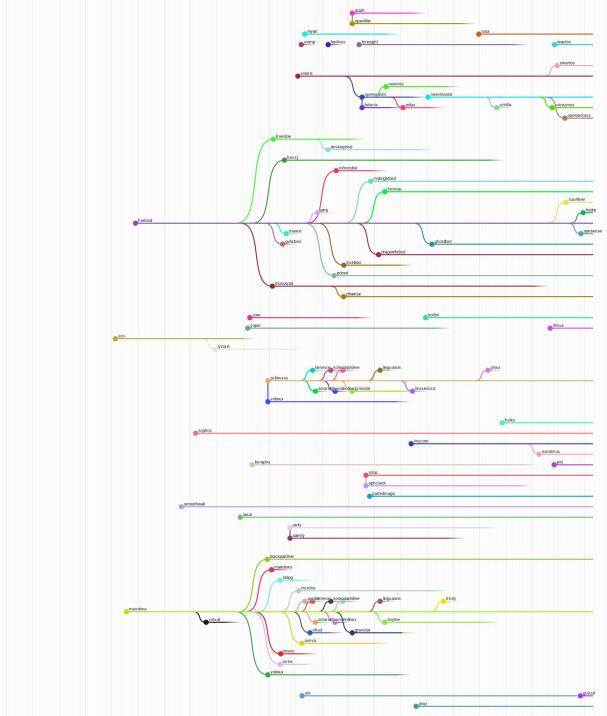
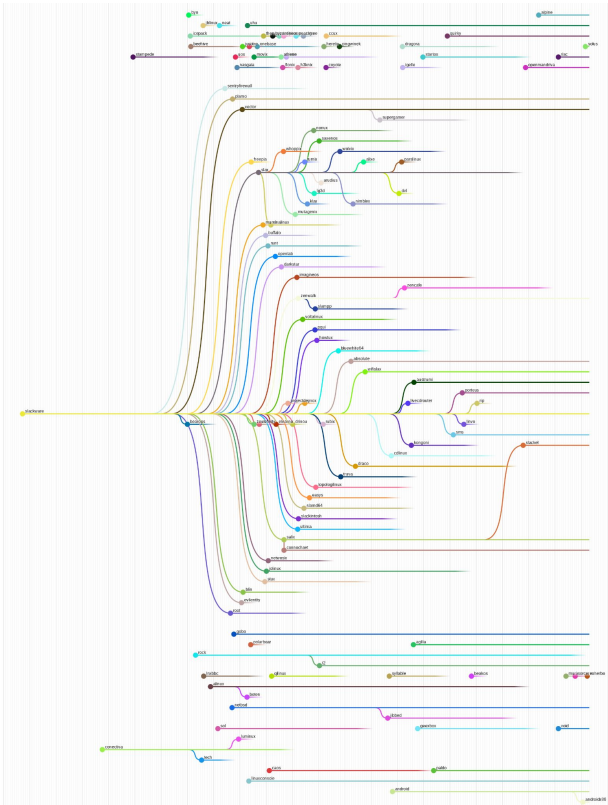
Michau/t5-base-en-generate-headline  
Text2Text Generation · Updated Jun 22 · ↓ 35k · ♥ 8



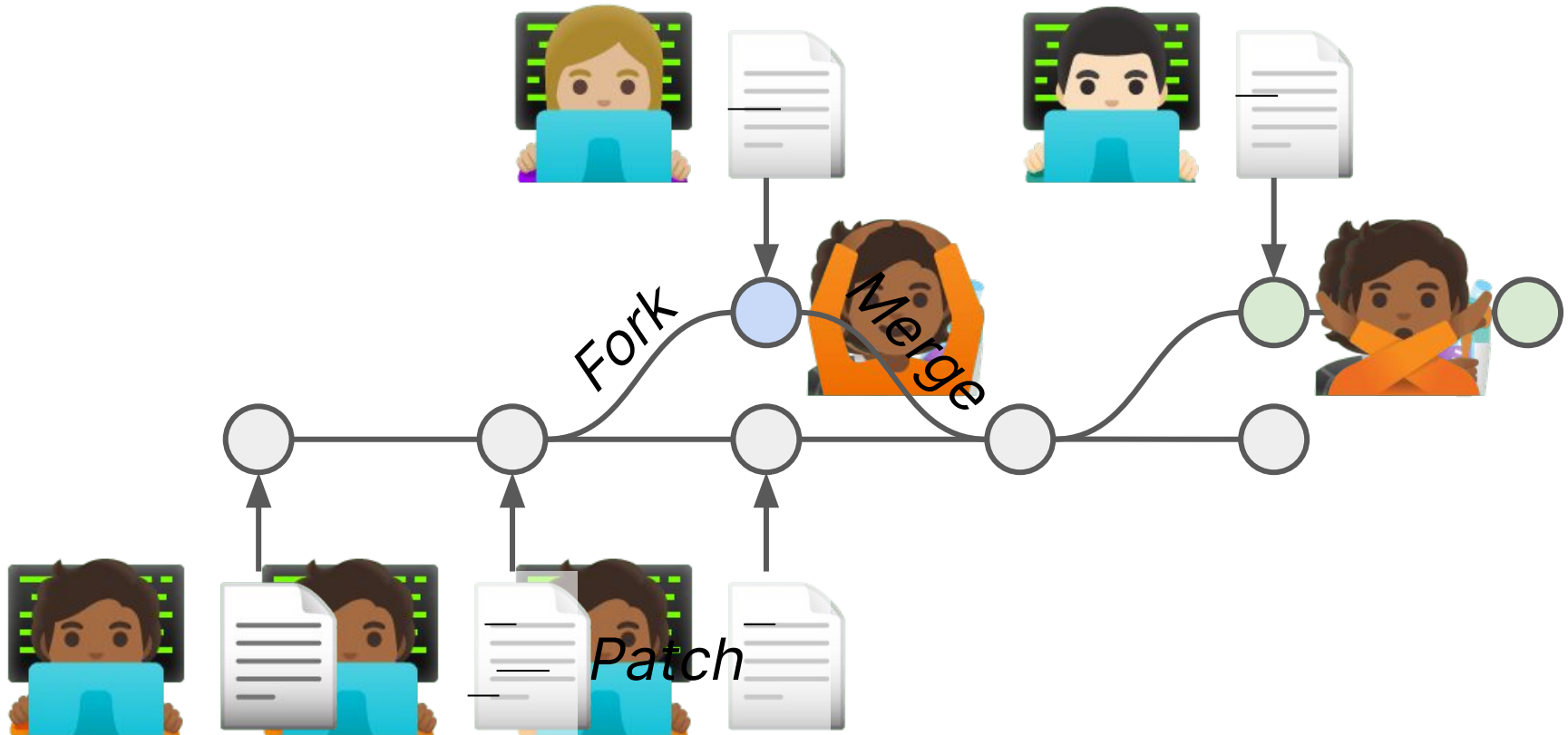
→ Additional training

- - - → New model

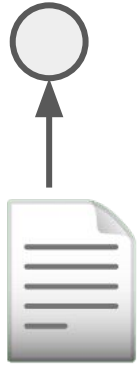


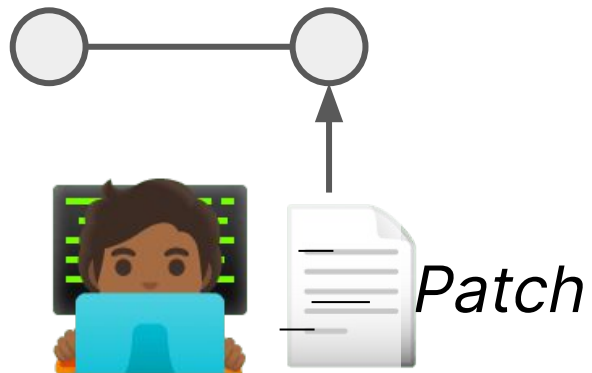


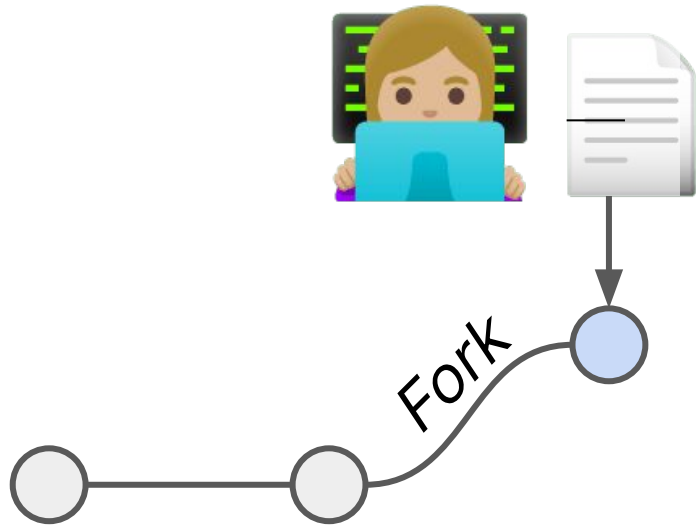
from <https://distrowatch.com/>

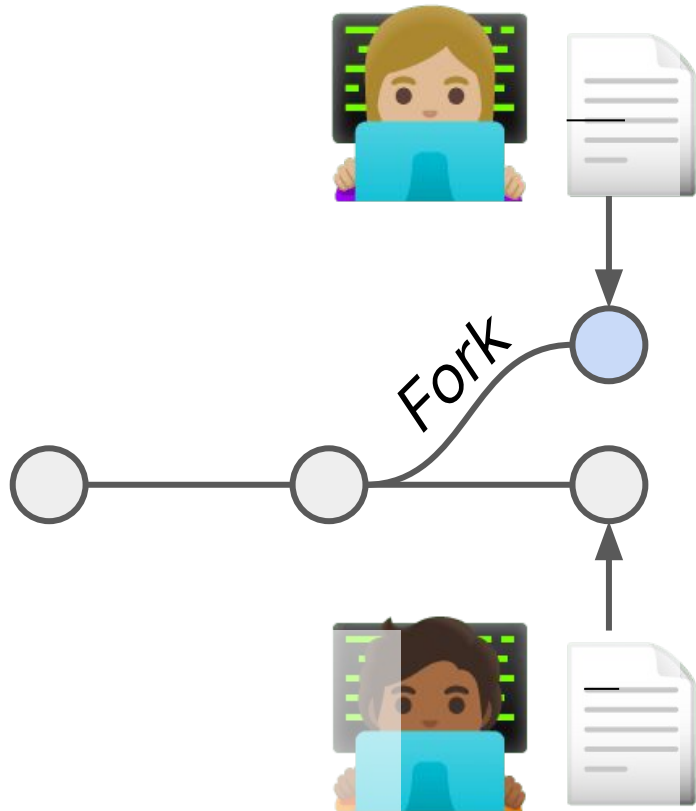


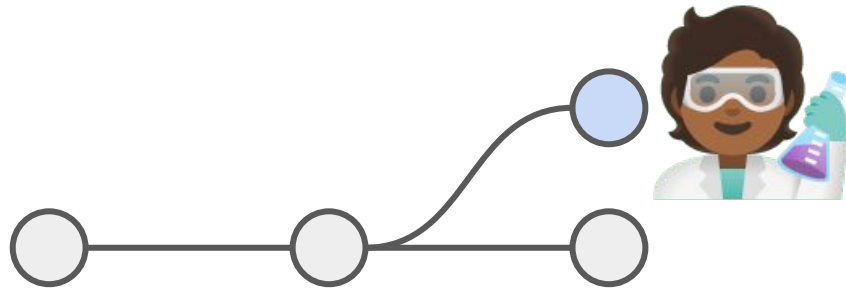


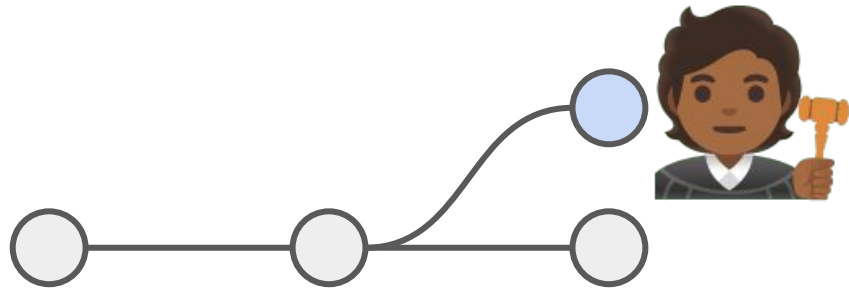


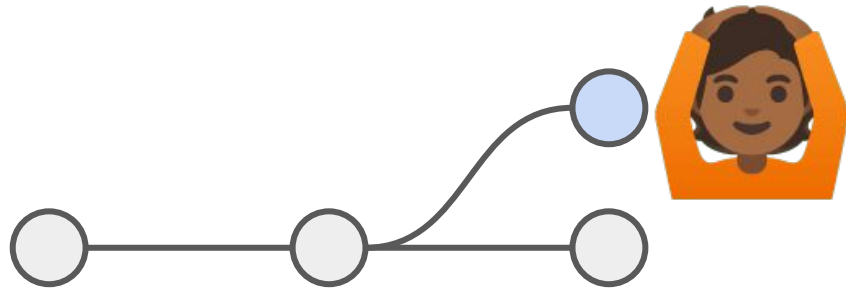


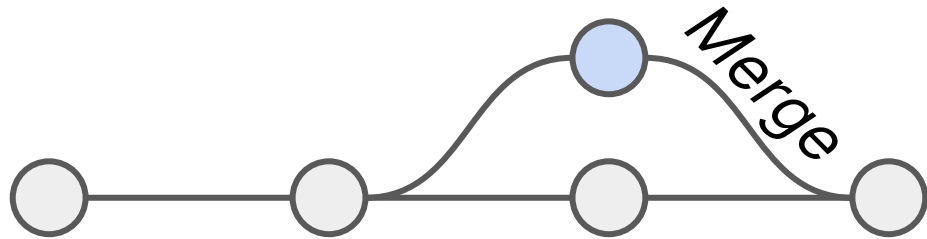




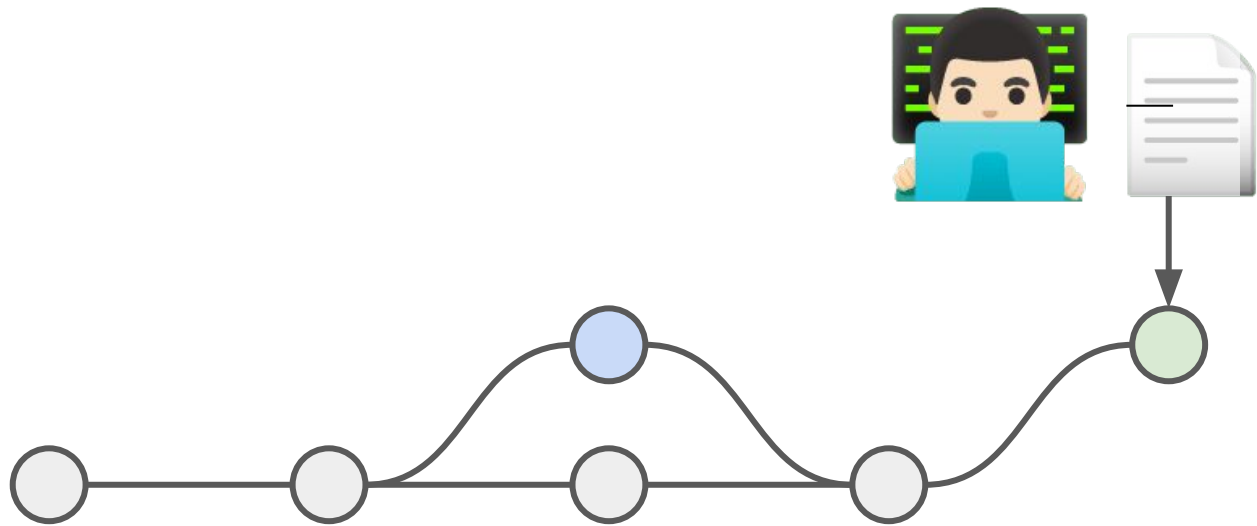


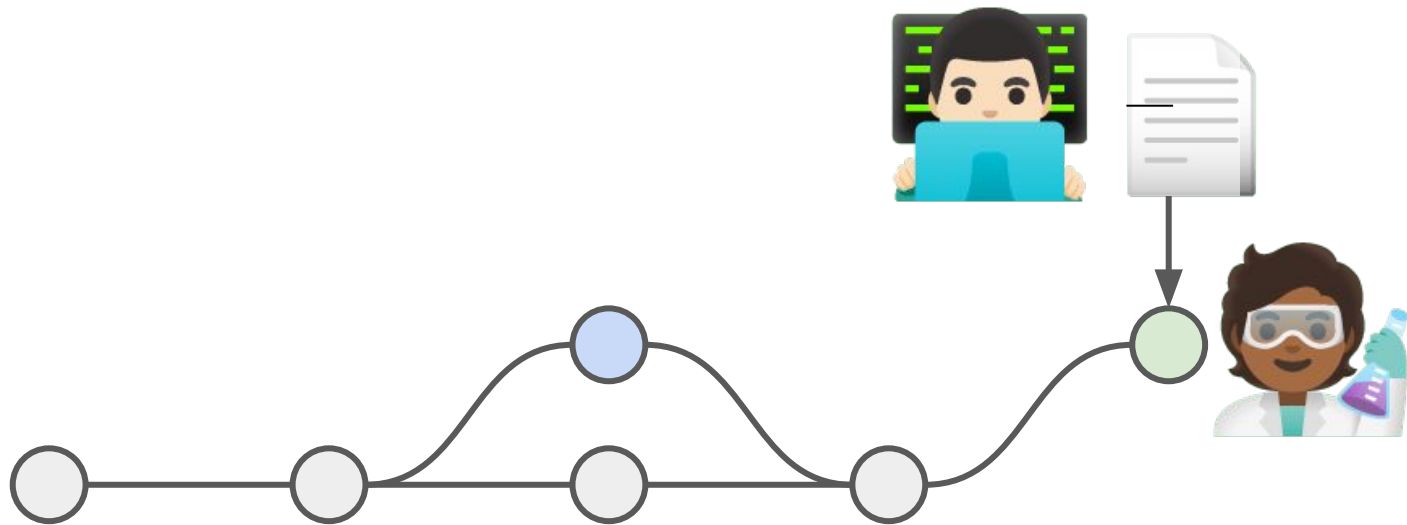


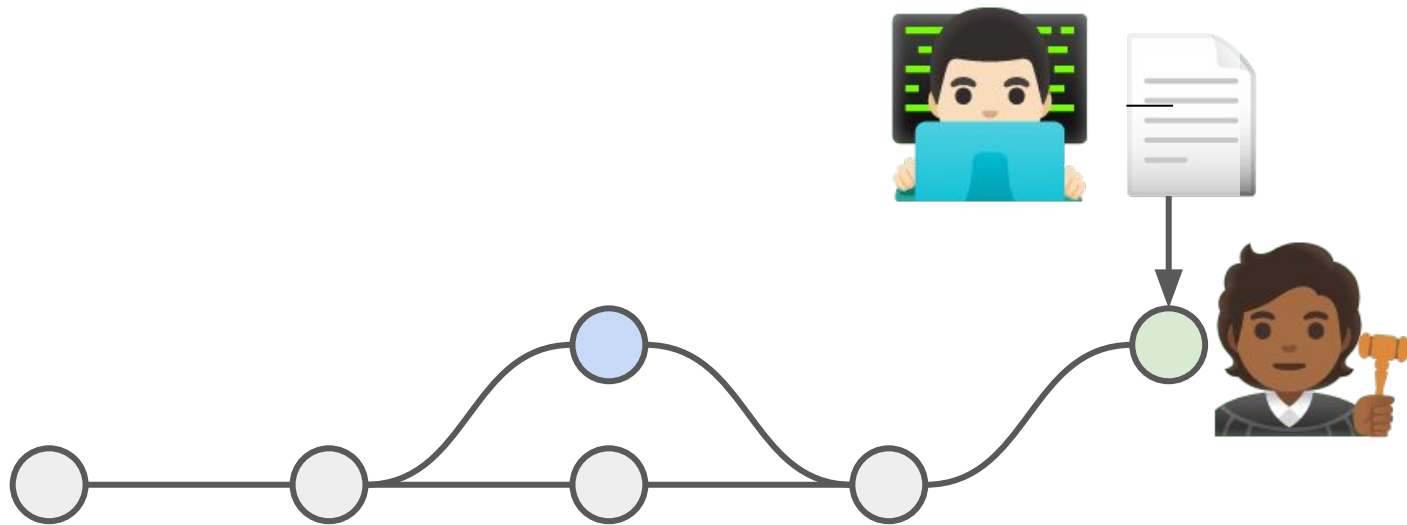


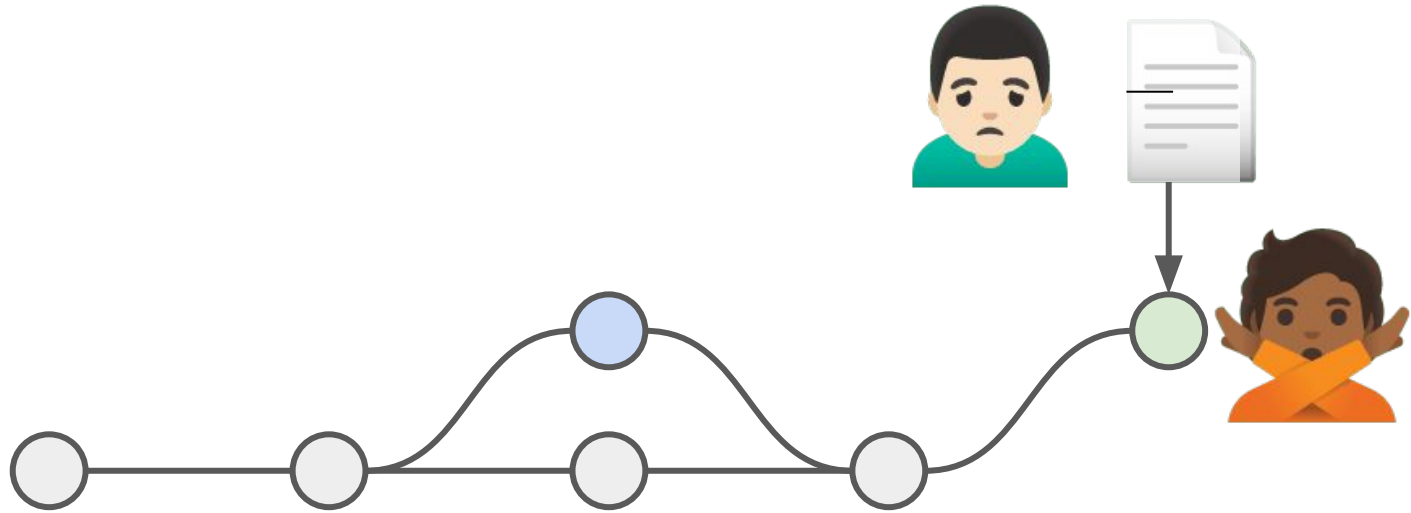


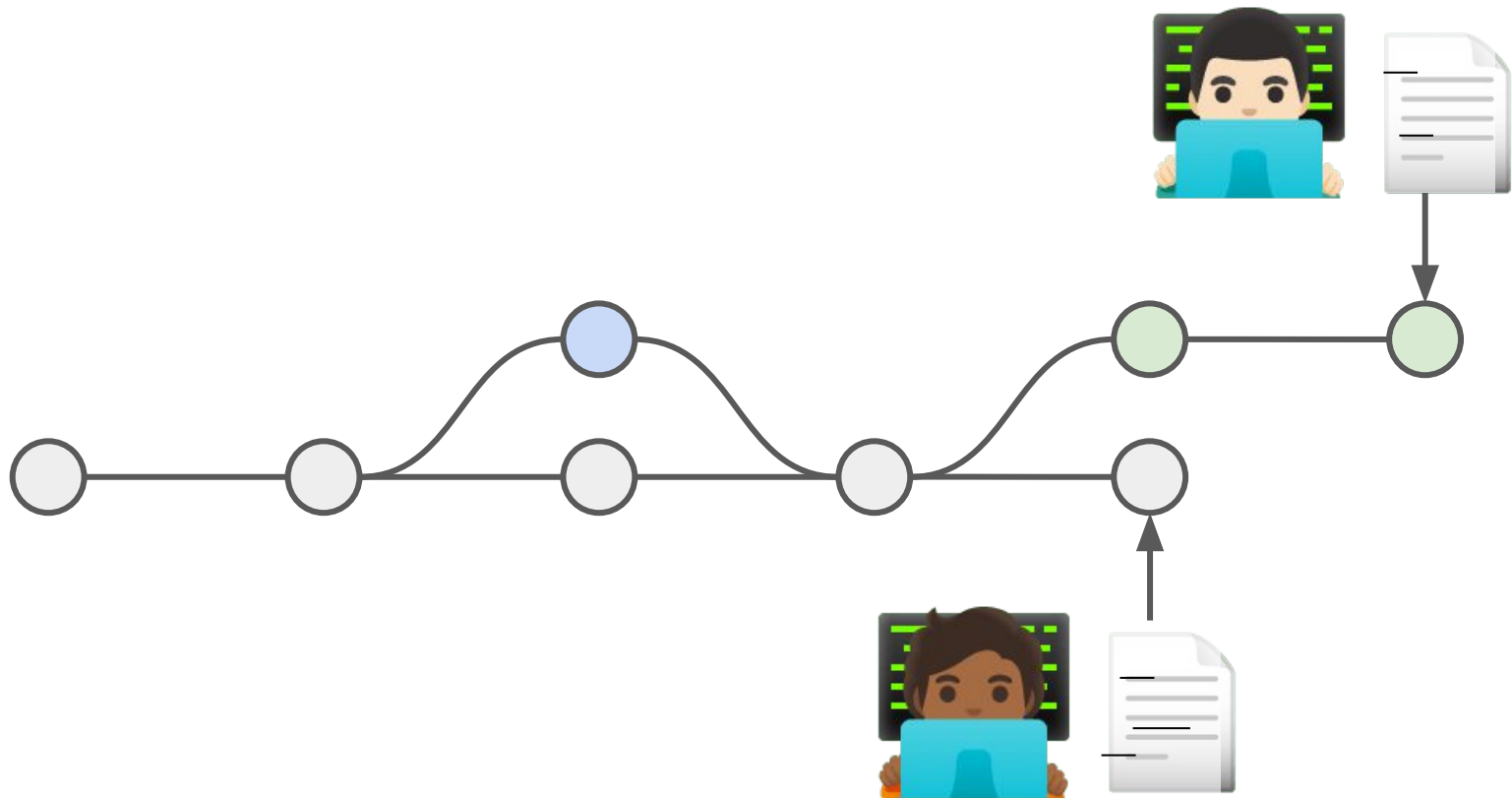


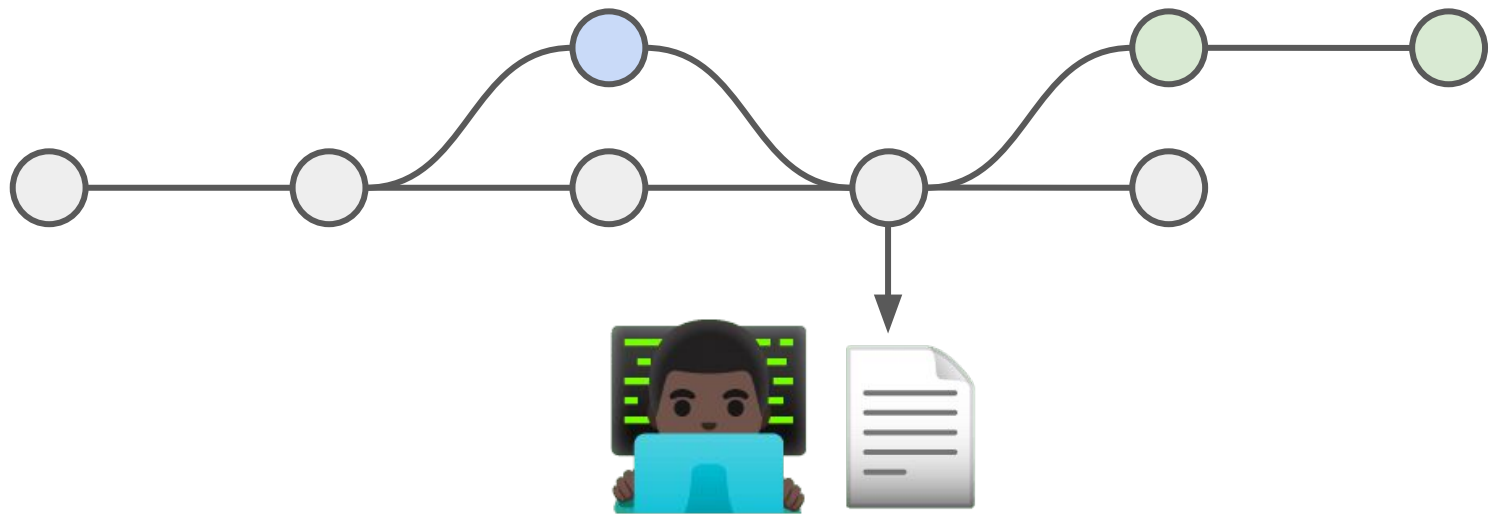












mlflow™



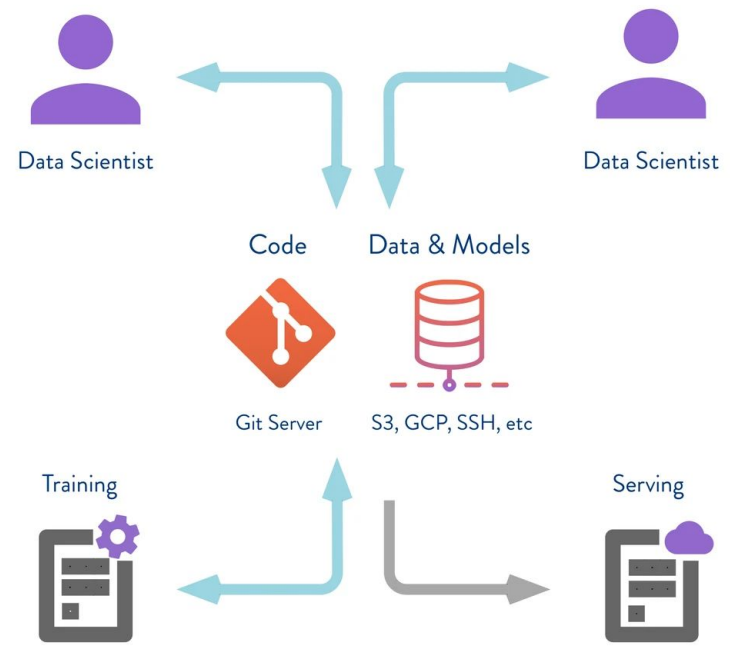
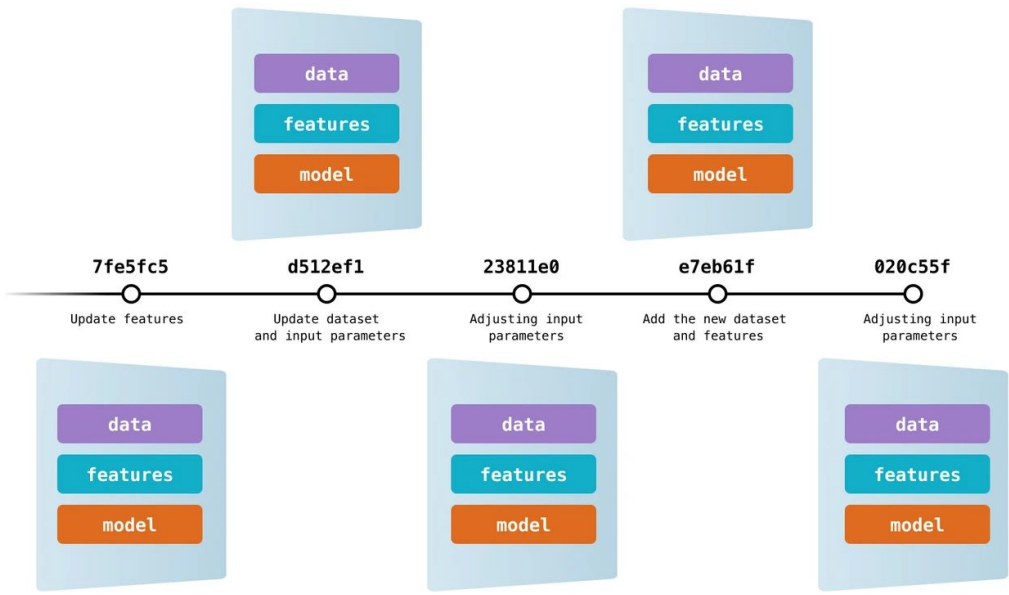
W&B



HUGGING FACE



comet



from <https://dvc.org/>

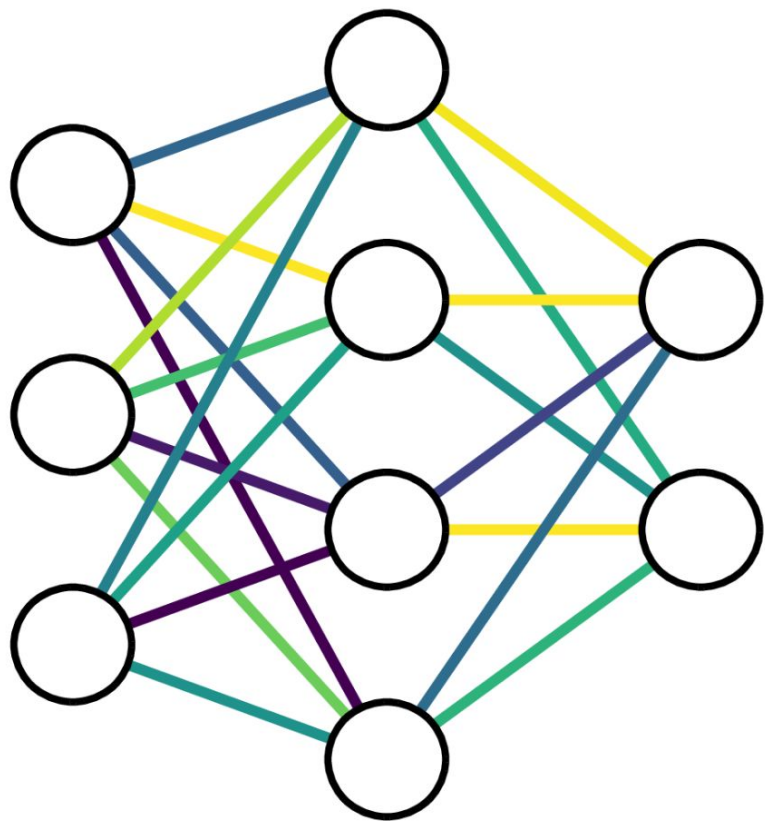
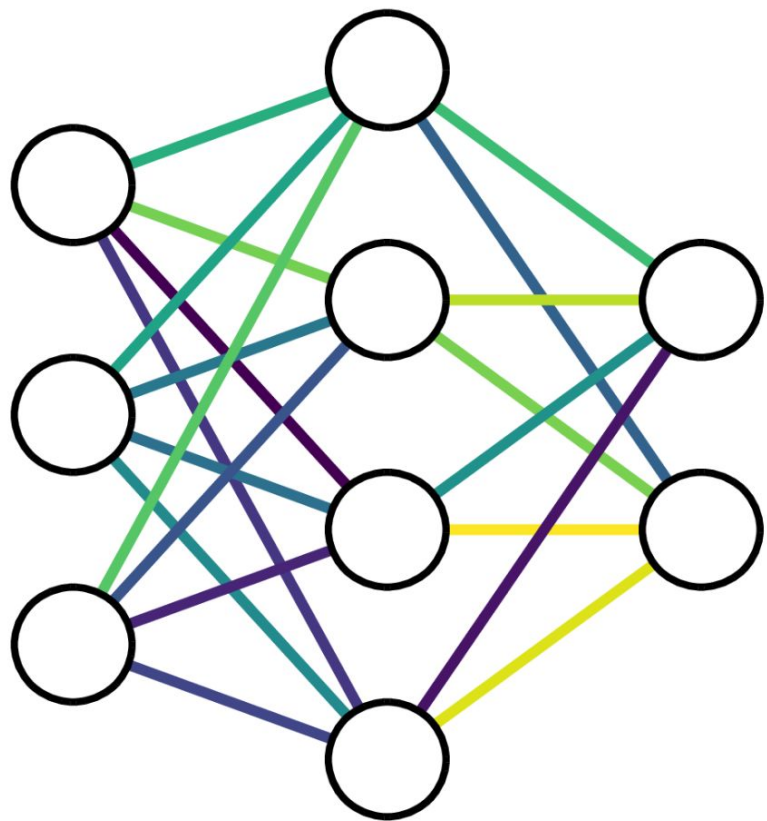


*How can we enable collaborative and continual development of machine learning models?*

We need to be able to cheaply communicate **patches** and **merge** updates from different contributors.

*How can we enable collaborative and continual development of machine learning models?*

We need to be able to cheaply communicate **patches** and **merge** updates from different contributors.





$$D_{\text{KL}}(p_{\theta}(y|x) \parallel p_{\theta+\delta}(y|x))$$

$$D_{\text{KL}}(p_{\theta}(y|x) \parallel p_{\theta+\delta}(y|x))$$

$$\mathbb{E}_x D_{\text{KL}}(p_{\theta}(y|x) \parallel p_{\theta+\delta}(y|x)) = \delta^{\text{T}} F_{\theta} \delta + O(\delta^3)$$

$$D_{\text{KL}}(p_{\theta}(y|x) \parallel p_{\theta+\delta}(y|x))$$

$$\mathbb{E}_x D_{\text{KL}}(p_{\theta}(y|x) \parallel p_{\theta+\delta}(y|x)) = \delta^{\text{T}} F_{\theta} \delta + O(\delta^3)$$

$$F_{\theta} = \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{y \sim p_{\theta}(y|x)} \nabla_{\theta} \log p_{\theta}(y|x) \nabla_{\theta} \log p_{\theta}(y|x)^{\text{T}} \right]$$

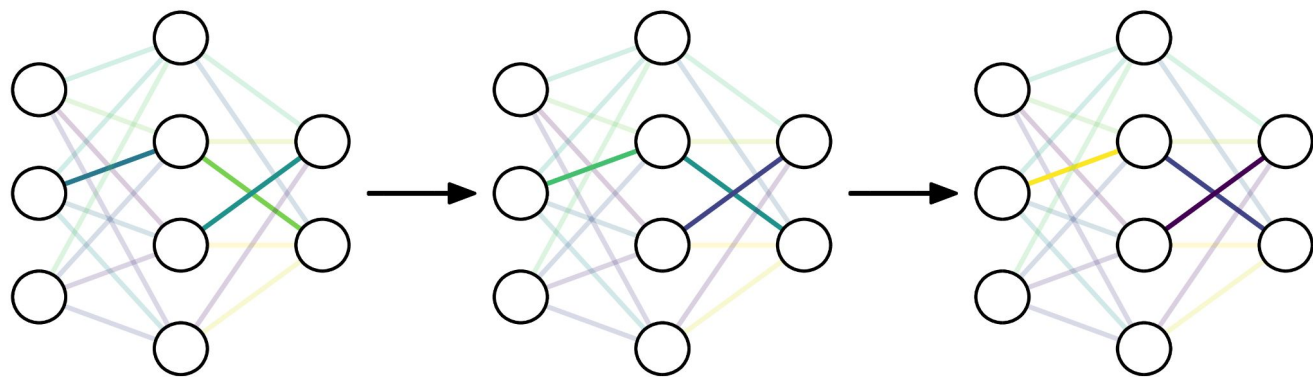
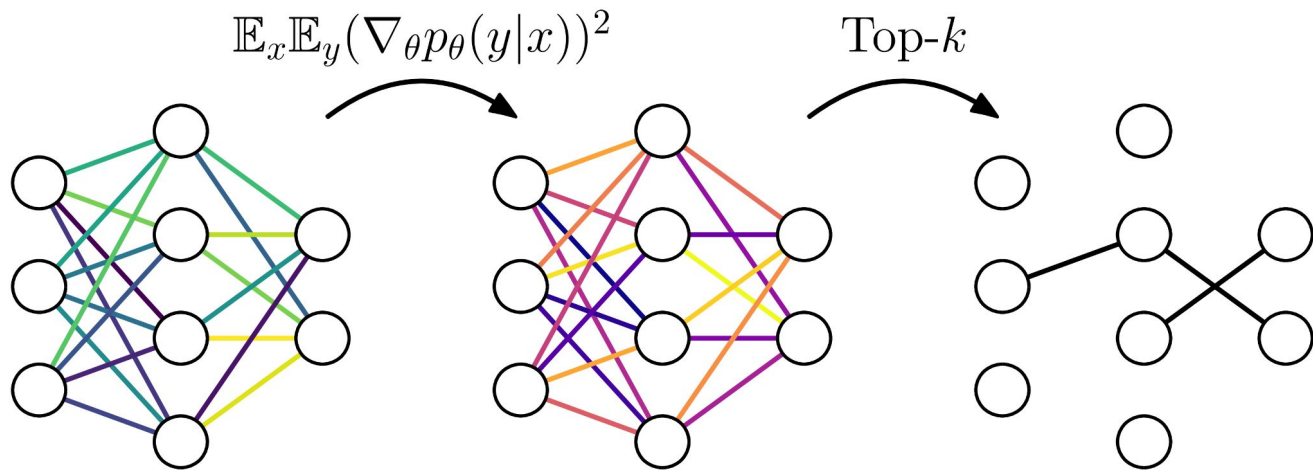
$$D_{\text{KL}}(p_{\theta}(y|x) \parallel p_{\theta+\delta}(y|x))$$

$$\mathbb{E}_x D_{\text{KL}}(p_{\theta}(y|x) \parallel p_{\theta+\delta}(y|x)) = \delta^{\text{T}} F_{\theta} \delta + O(\delta^3)$$

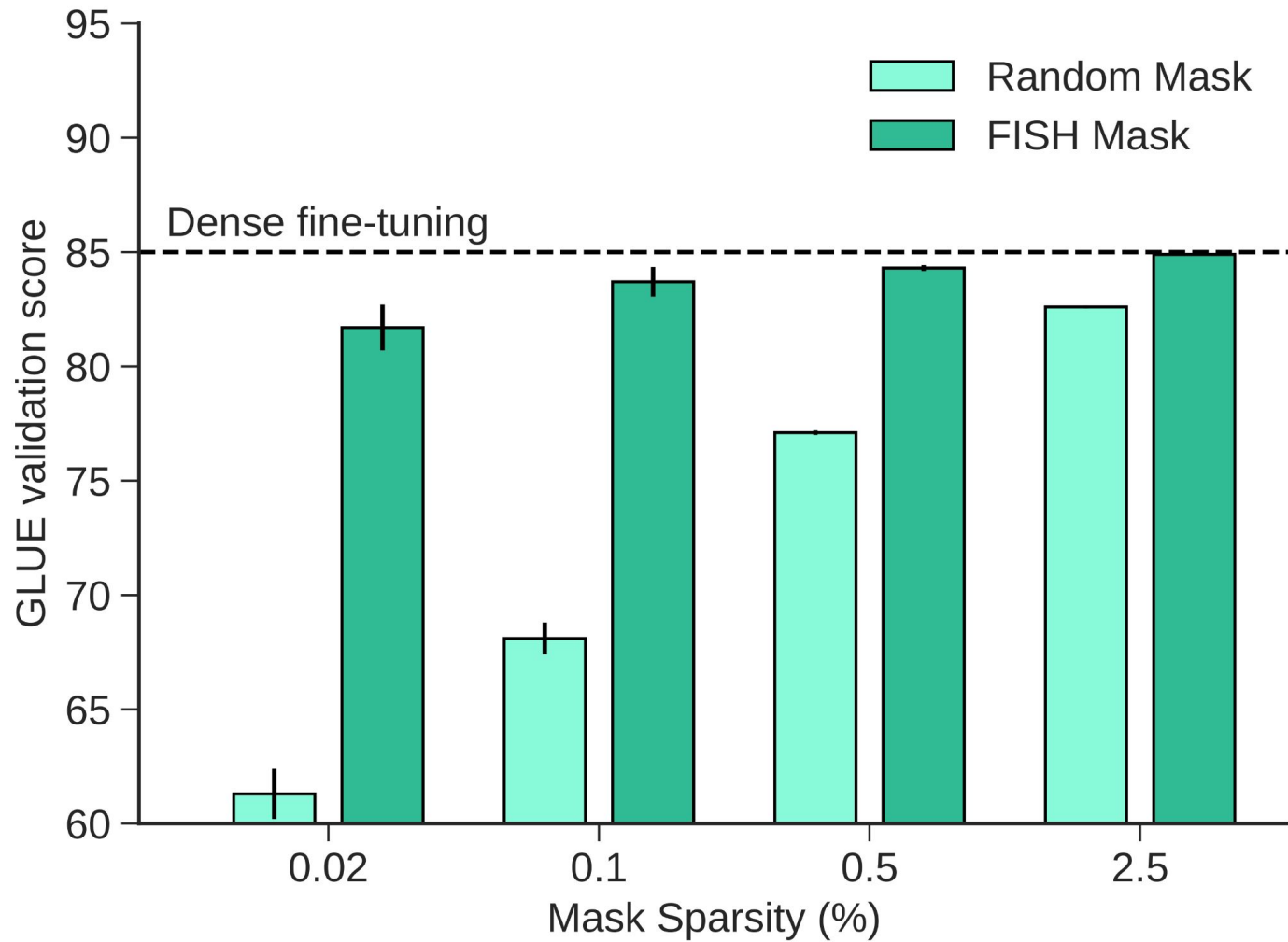
$$F_{\theta} = \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{y \sim p_{\theta}(y|x)} \nabla_{\theta} \log p_{\theta}(y|x) \nabla_{\theta} \log p_{\theta}(y|x)^{\text{T}} \right]$$

$$\hat{F}_{\theta} = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{y \sim p_{\theta}(y|x_i)} (\nabla_{\theta} \log p_{\theta}(y|x_i))^2$$





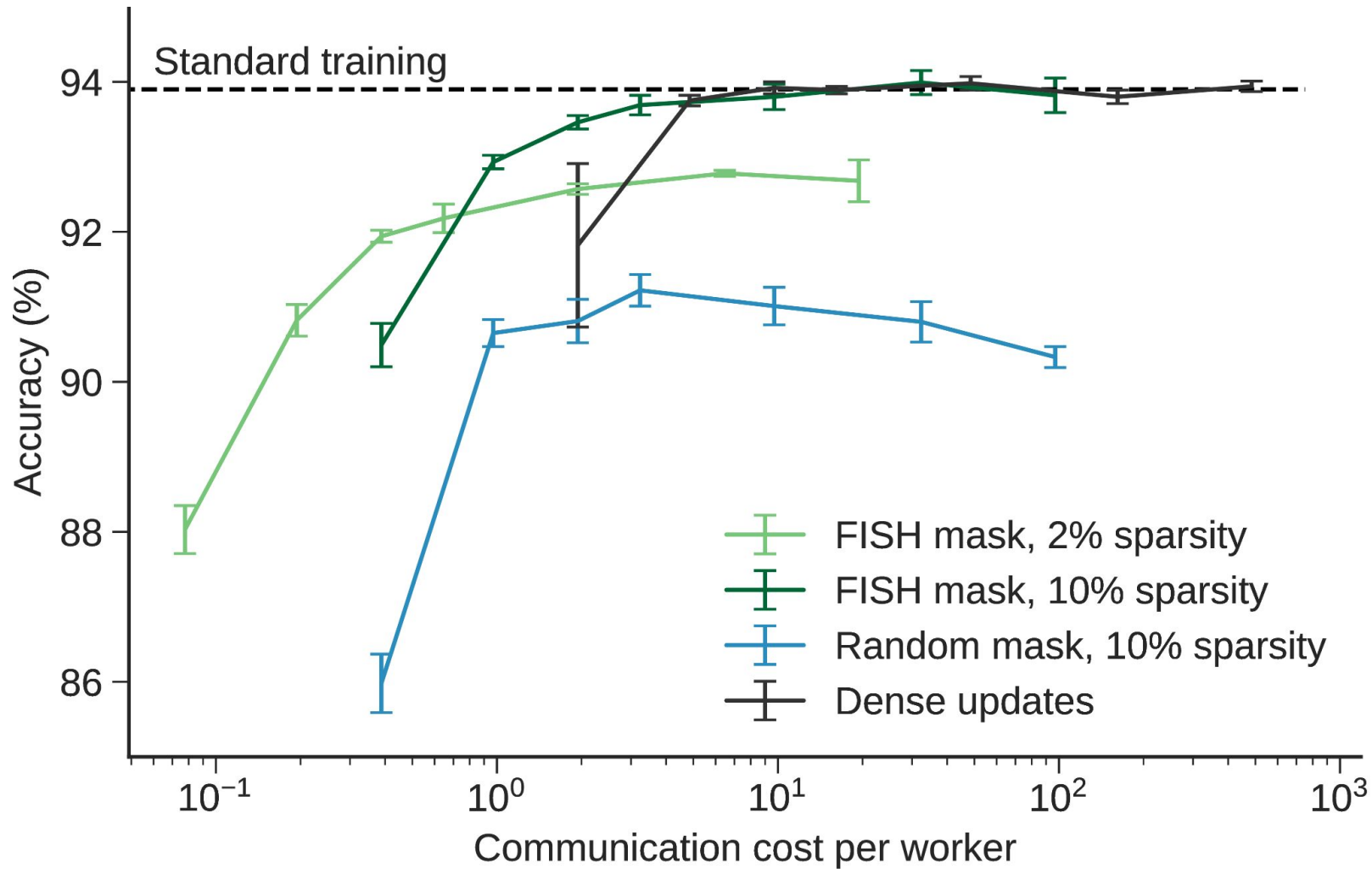
**Fisher-Induced Sparse Unchanging (FISH) Mask**



---

Method	Sparsity	GLUE Score
Dense Fine-tuning	100%	82.5
Bit-Fit	0.08%	81.2
FISH Mask	0.08%	81.3
Diff Pruning	0.50%	81.5
FISH Mask	0.50%	82.6

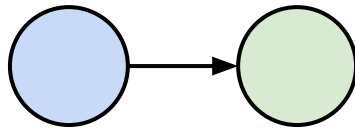
---



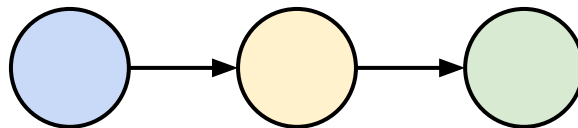
*How can we enable collaborative and continual development of machine learning models?*

We need to be able to cheaply communicate **patches** and **merge** updates from different contributors.

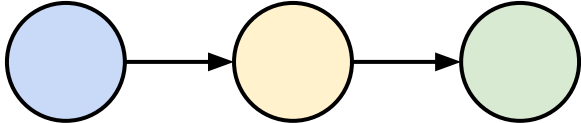
Pre-training    Downstream



Pre-training    Intermediate    Downstream

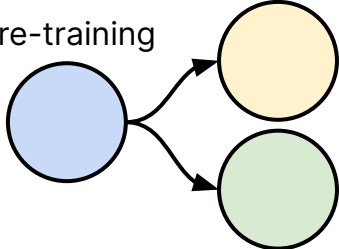


Pre-training    Intermediate    Downstream



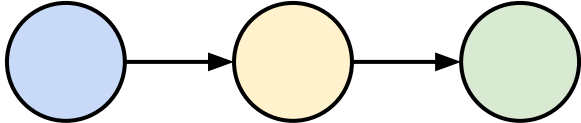
Intermediate

Pre-training



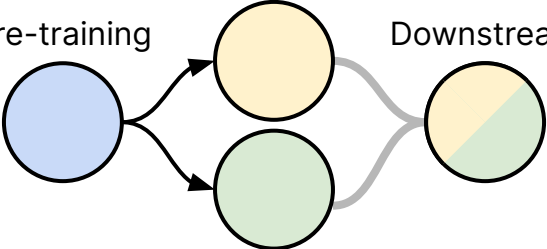


Pre-training    Intermediate    Downstream

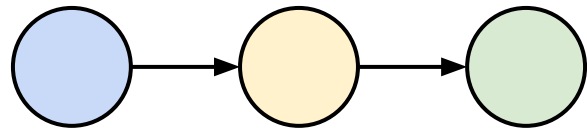


Intermediate

Pre-training    Downstream

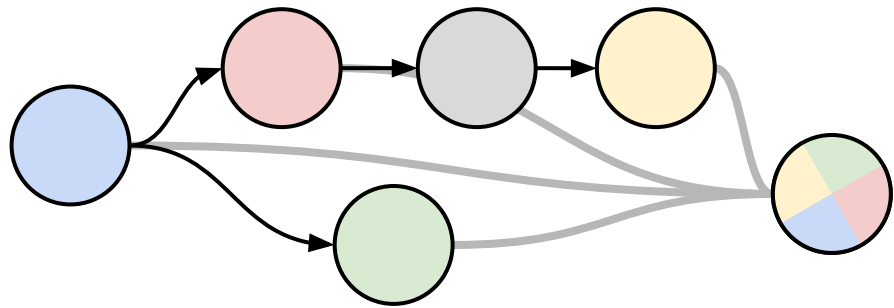
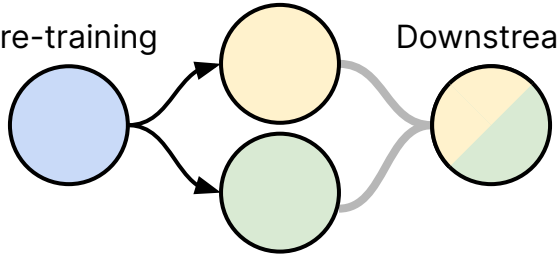


Pre-training    Intermediate    Downstream



Intermediate

Pre-training    Downstream



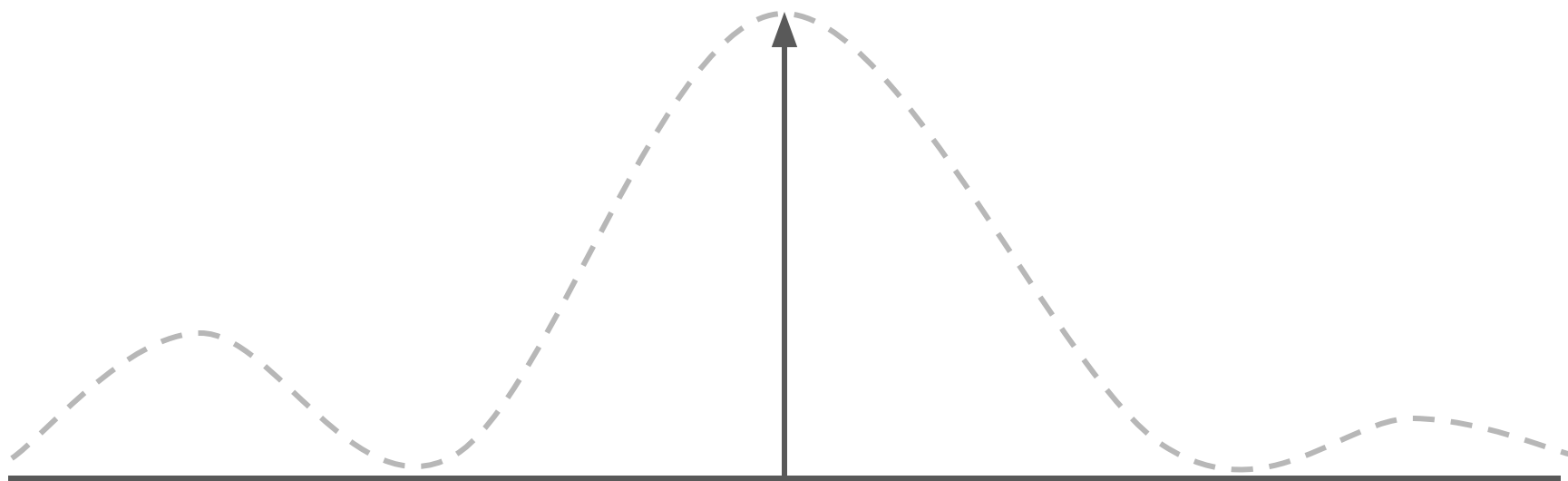
$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log p(\theta | \mathcal{D}_i)$$

$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \underbrace{\log p(\theta | \mathcal{D}_i)}_{\substack{\text{Log posterior} \\ \text{for model } i}}$$

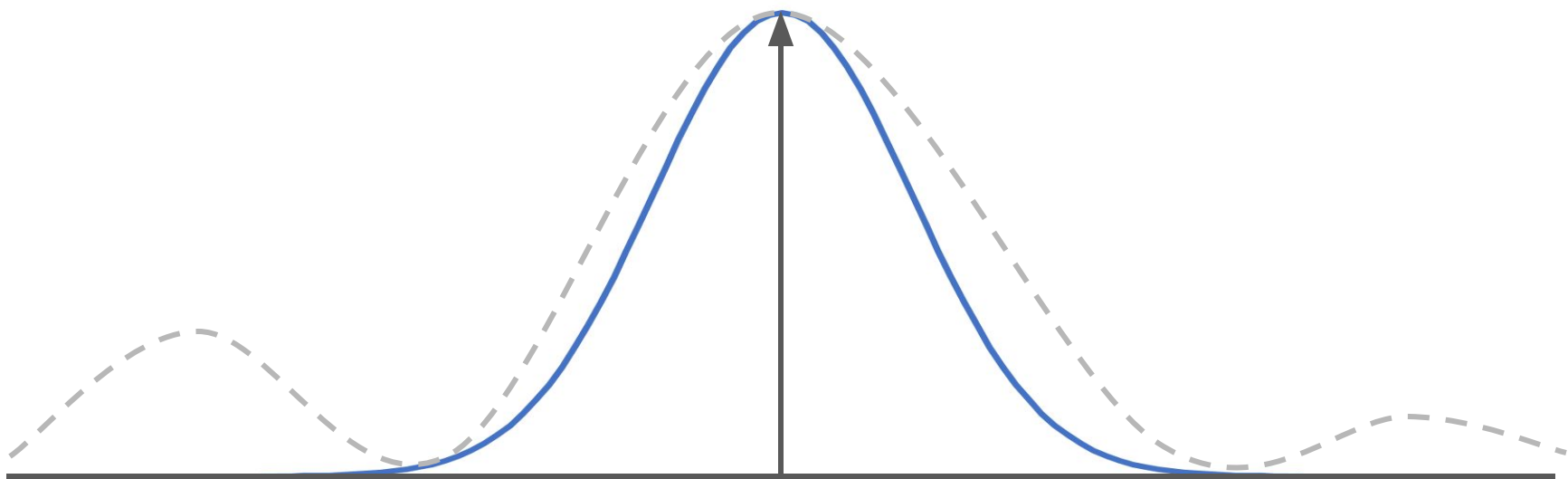
$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log p(\theta | \mathcal{D}_i)$$

*Hyperparameter  
controlling the  
importance of model  $i$*

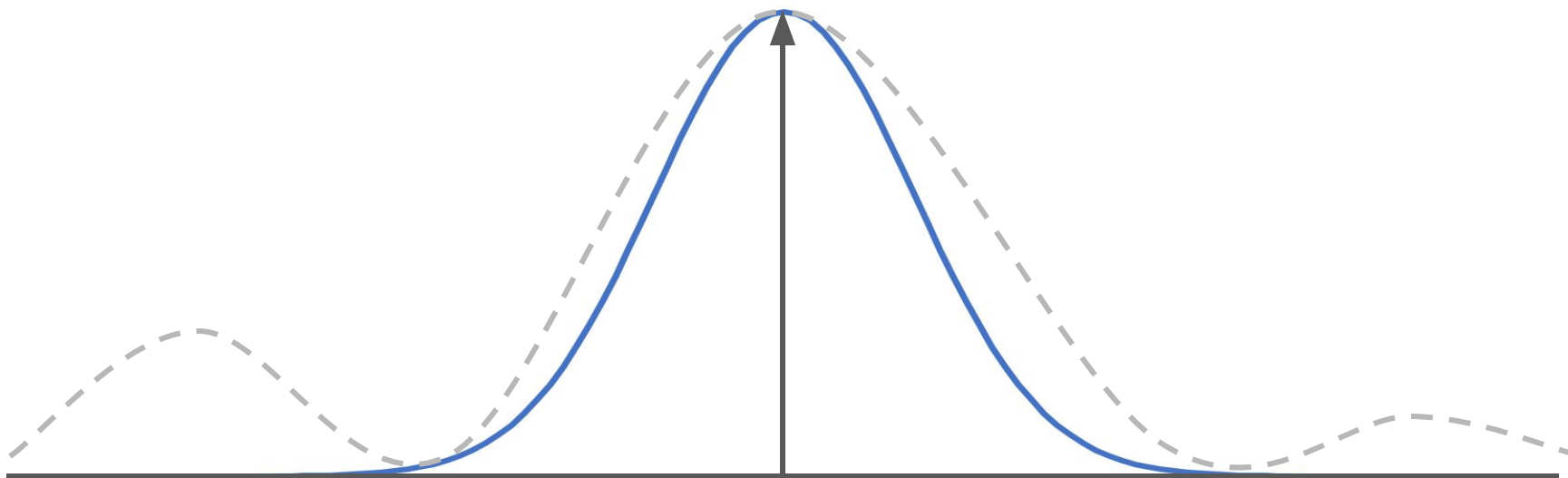
$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log p(\text{👤})$$



$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log \mathcal{N}(\theta | \theta_i, H_i^{-1})$$

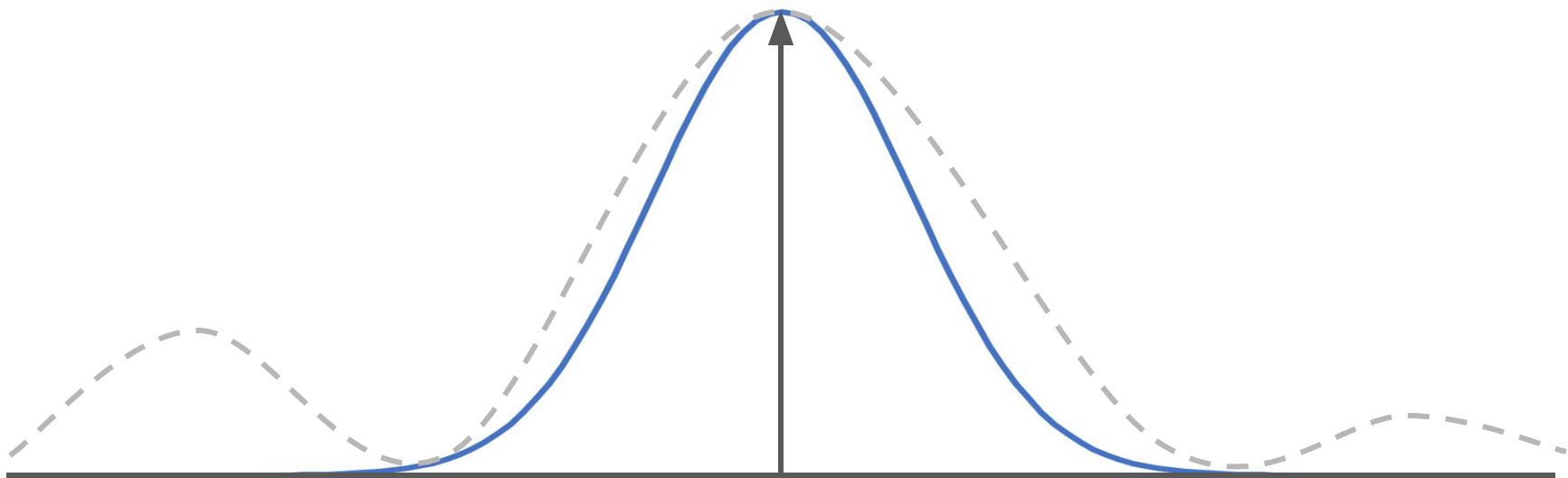


$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log \mathcal{N}(\theta | \theta_i, \sigma^2)$$

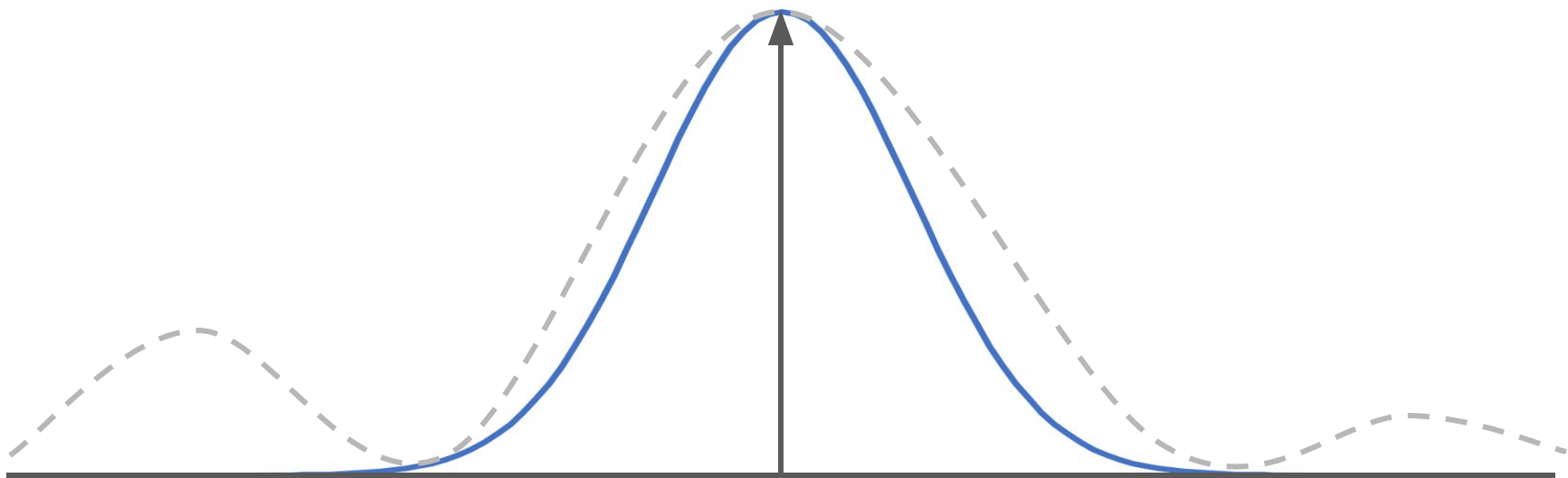




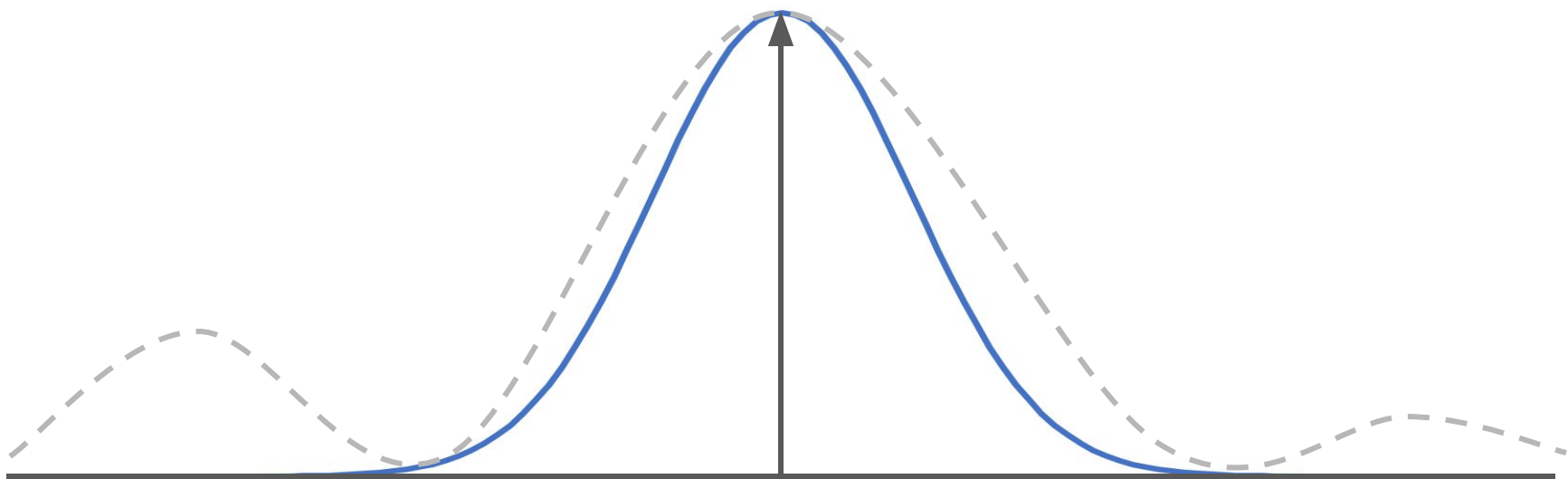
$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log \mathcal{N}(\theta | \theta_i, F_i)$$



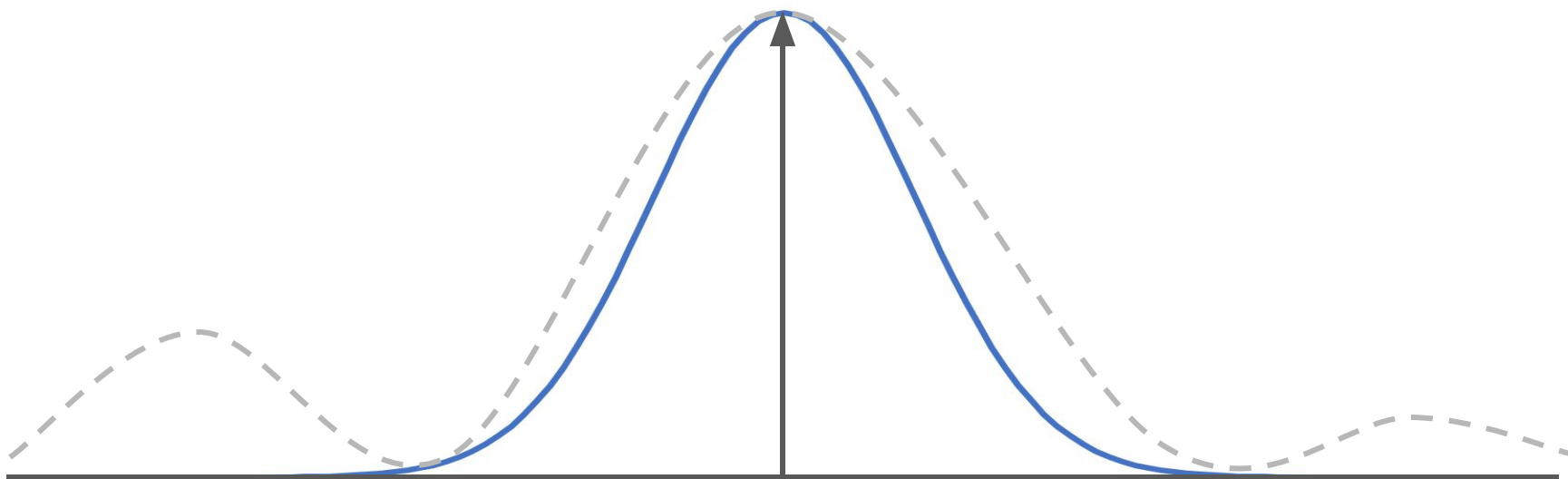
$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log \mathcal{N}(\theta | \theta_i)$$

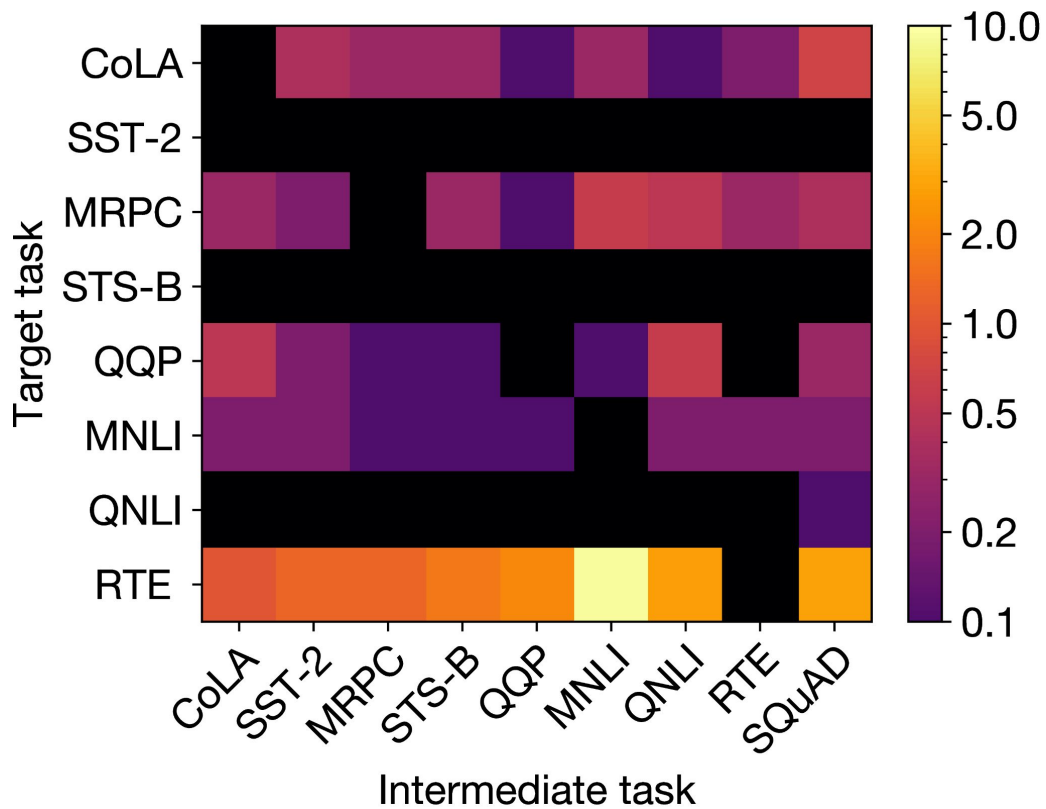
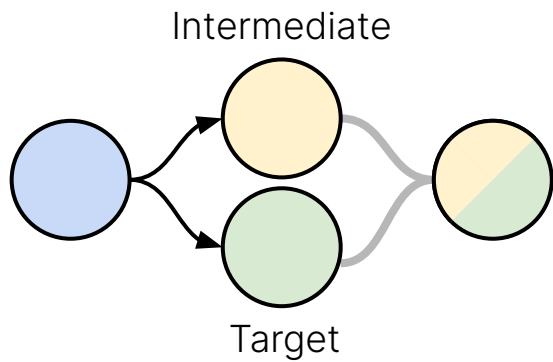


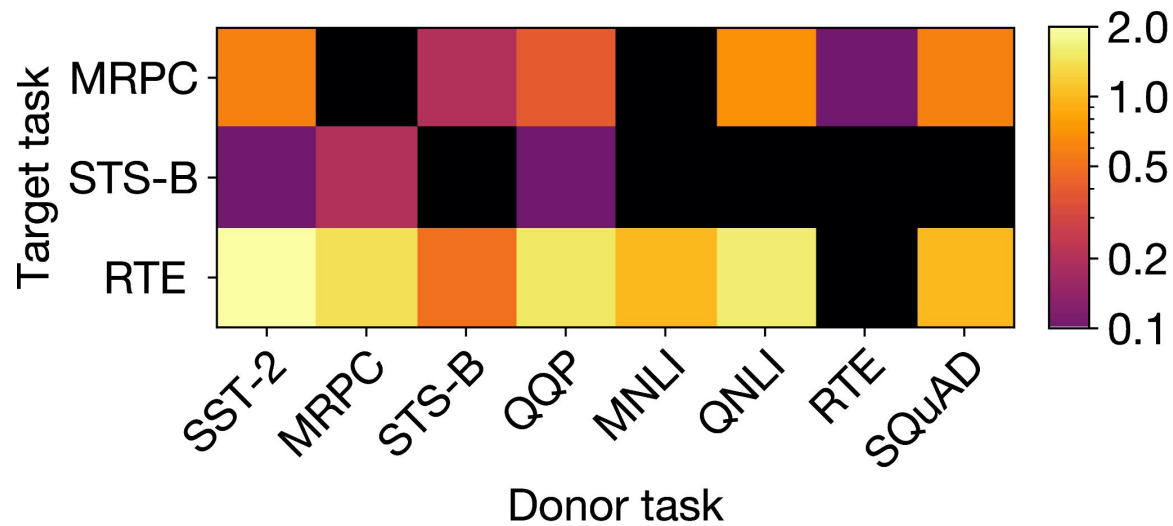
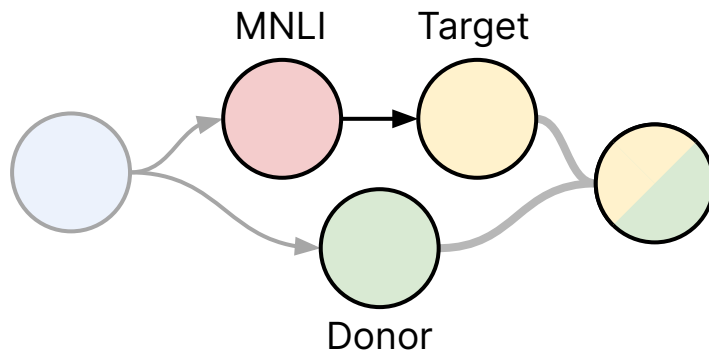
$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log \mathcal{N}(\theta | \theta_i, \hat{F}_i)$$

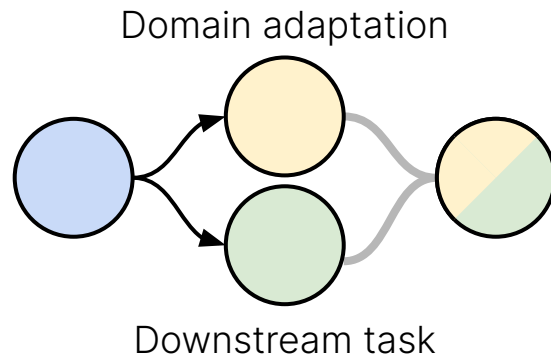


$$\arg \max_{\theta} \sum_{i=1}^M \lambda_i \log \mathcal{N}(\theta | \theta_i, \hat{F}_i)$$









Task	Unmerged	Merged	Fine-tuned
CHEMPROT	82.7 <sub>0.3</sub>	83.1 <sub>0.4</sub>	82.5 <sub>0.1</sub>
ACL-ARC	70.5 <sub>3.2</sub>	73.2 <sub>1.7</sub>	71.5 <sub>3.0</sub>
SCIERC	81.0 <sub>0.4</sub>	81.3 <sub>0.5</sub>	81.6 <sub>1.0</sub>

*How can we enable collaborative and continual development of machine learning models?*

We need to be able to cheaply communicate **patches** and **merge** updates from different contributors.



*How can we enable collaborative and continual development of machine learning models?*

We need to be able to **rapidly evaluate** proposed changes to the model to ensure backward compatibility.

*How can we enable collaborative and continual development of machine learning models?*

We need to be able to combine **modular** components of different models to provide new skills and capabilities.

## Training Neural Networks with Fixed Sparse Masks

Yi-Lin Sung, Varun Nair, and Colin Raffel

## Merging Models with Fisher-Weighted Averaging

Michael Matena and Colin Raffel

Please give me feedback:

<http://bit.ly/colin-talk-feedback>