

# Meta-Learning

DSAA 2021 tutorial

Joaquin Vanschoren (TU Eindhoven)

# Intro: humans can easily learn from a single example

*thanks to years of learning (and eons of evolution)*

Canna Indica ‘Picasso’



**train**

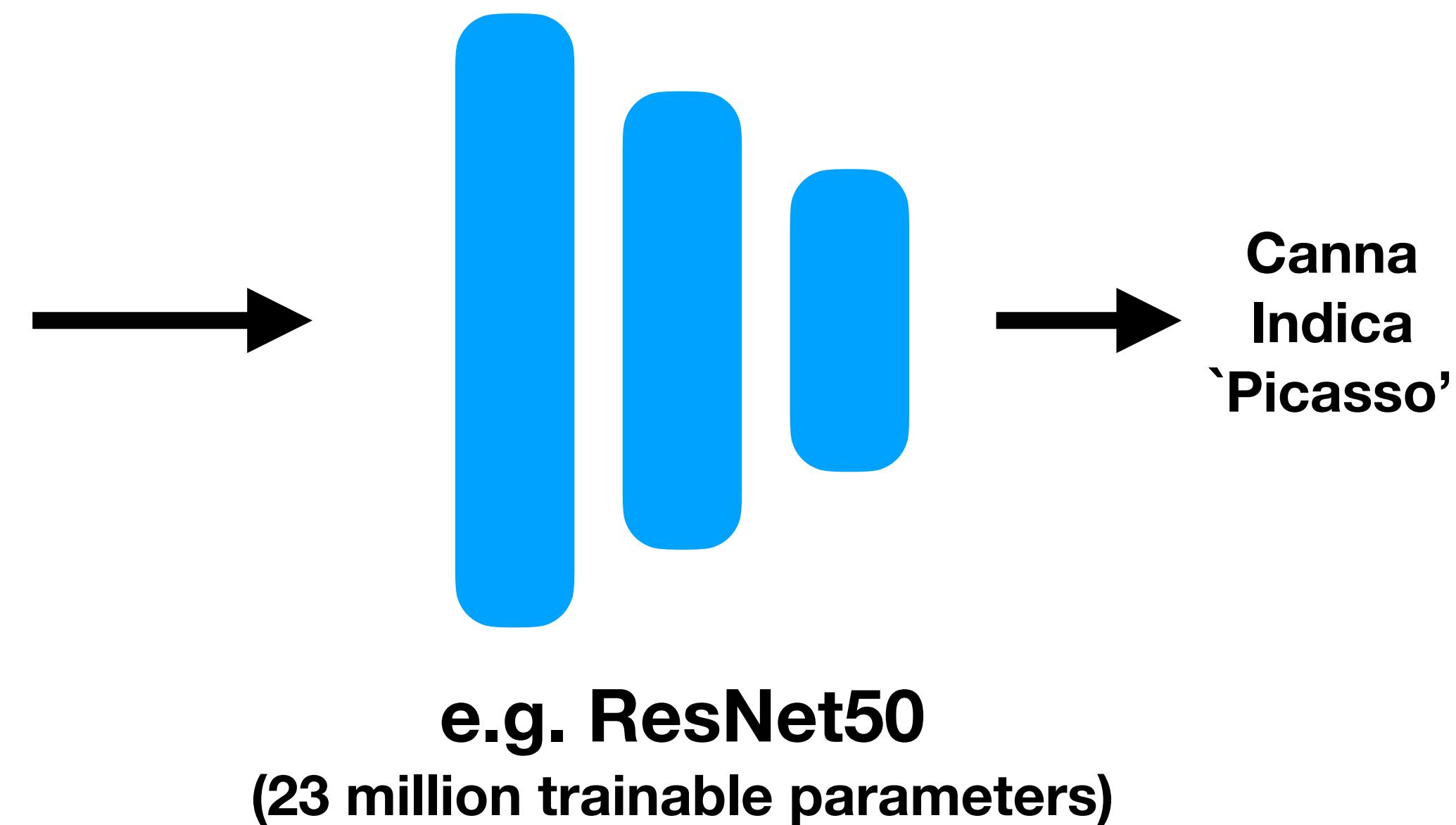
later  
→



**test**

?

# Can a computer learn from a single example?



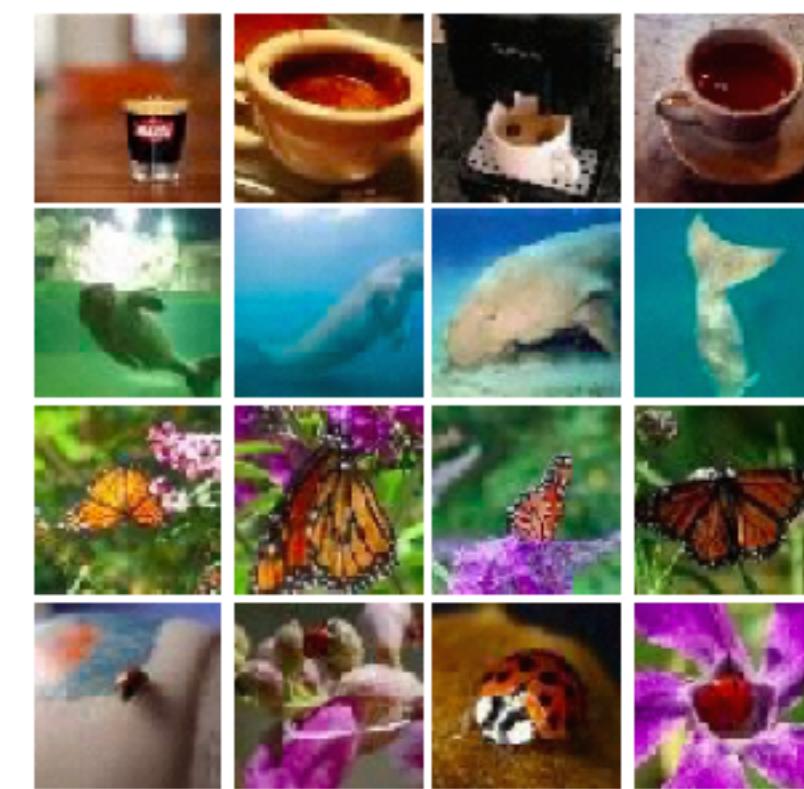
That won't work :) Humans also don't start from scratch.

# Transfer learning?

Target task

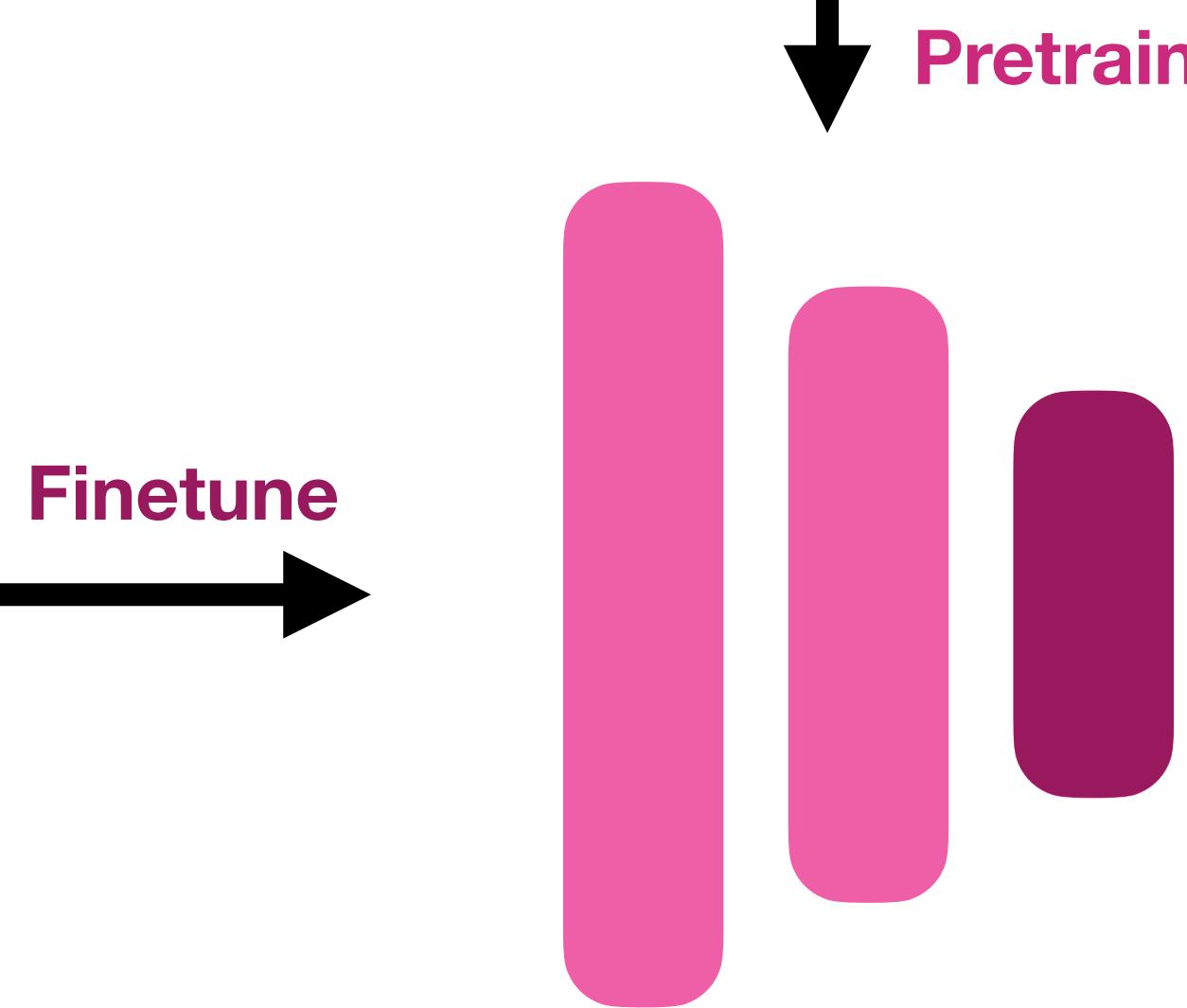


Source task



**ImageNet**  
(14 million images)

Finetune



e.g. ResNet50  
(23 million trainable parameters)

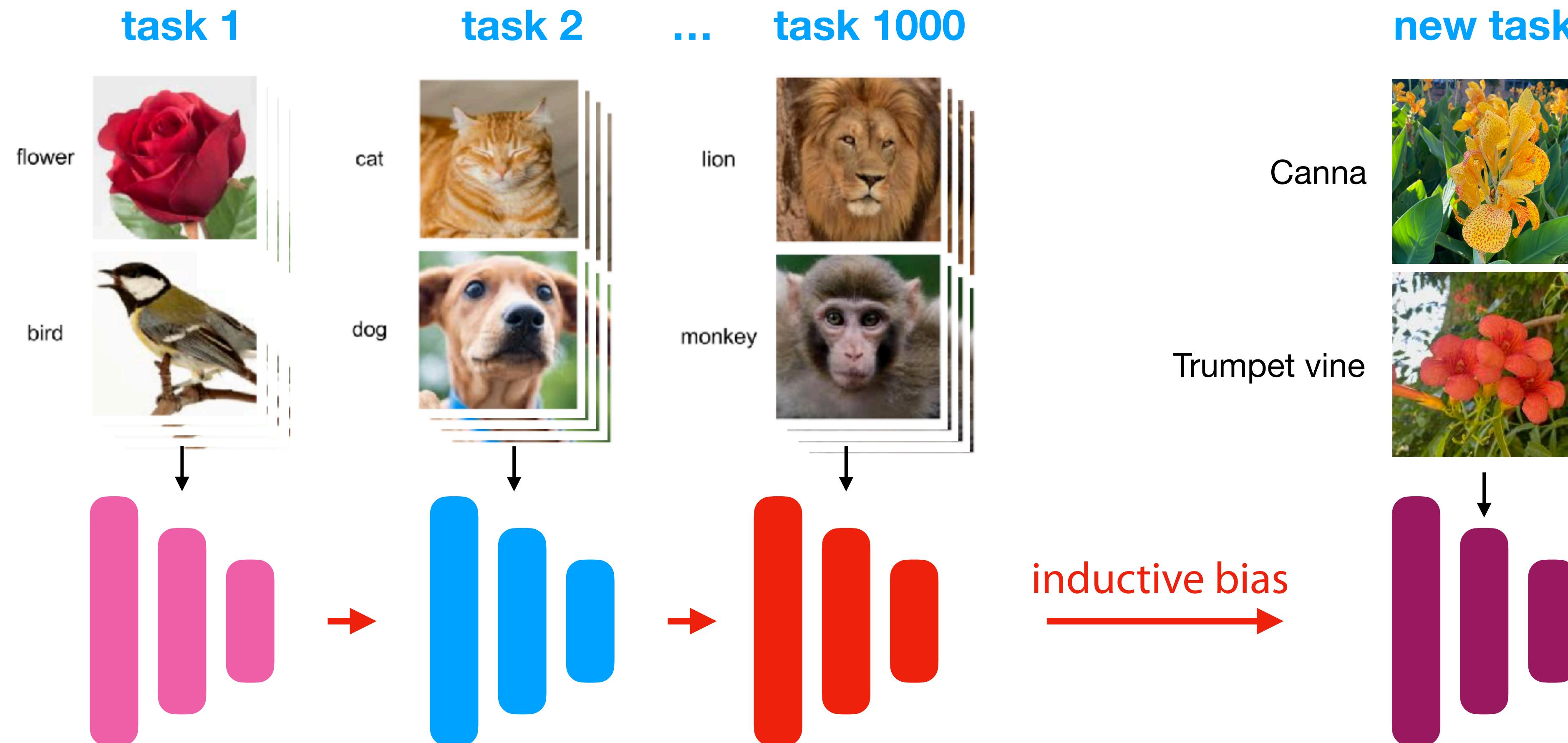
A single source task (e.g. ImageNet) may not generalize well to the test task

# Meta-learning

Learn over a series (or distribution) of many different tasks/episodes

**Inductive bias:** learn **assumptions** that you can transfer to new tasks

Prepare yourself to learn new things faster



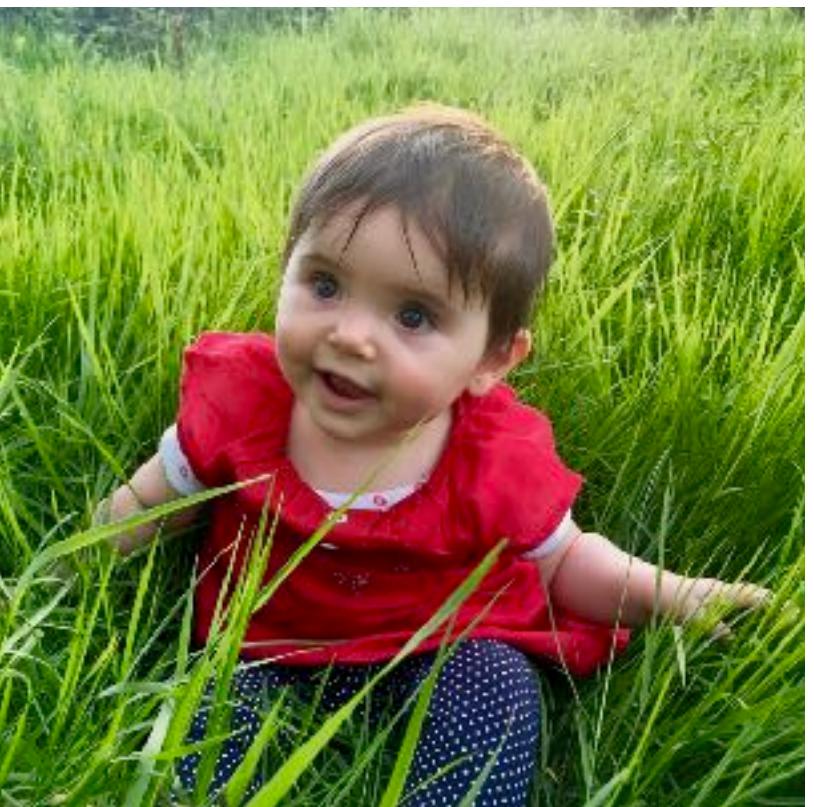
Useful in many real-life situations: rare events, test-time constraints, data collection costs, privacy issues,...

# Inspired by human learning

We don't transfer from a single source task, we learn across many, many tasks

We have a 'drive' to explore new, challenging, but doable, fun tasks

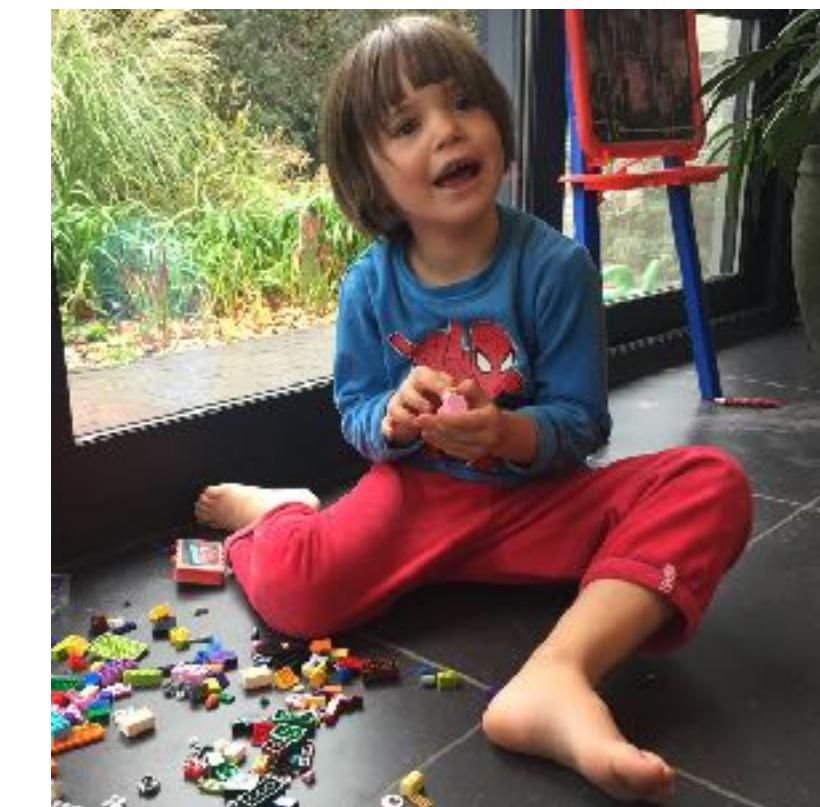
**year 1**



**year 2**



**year 3**

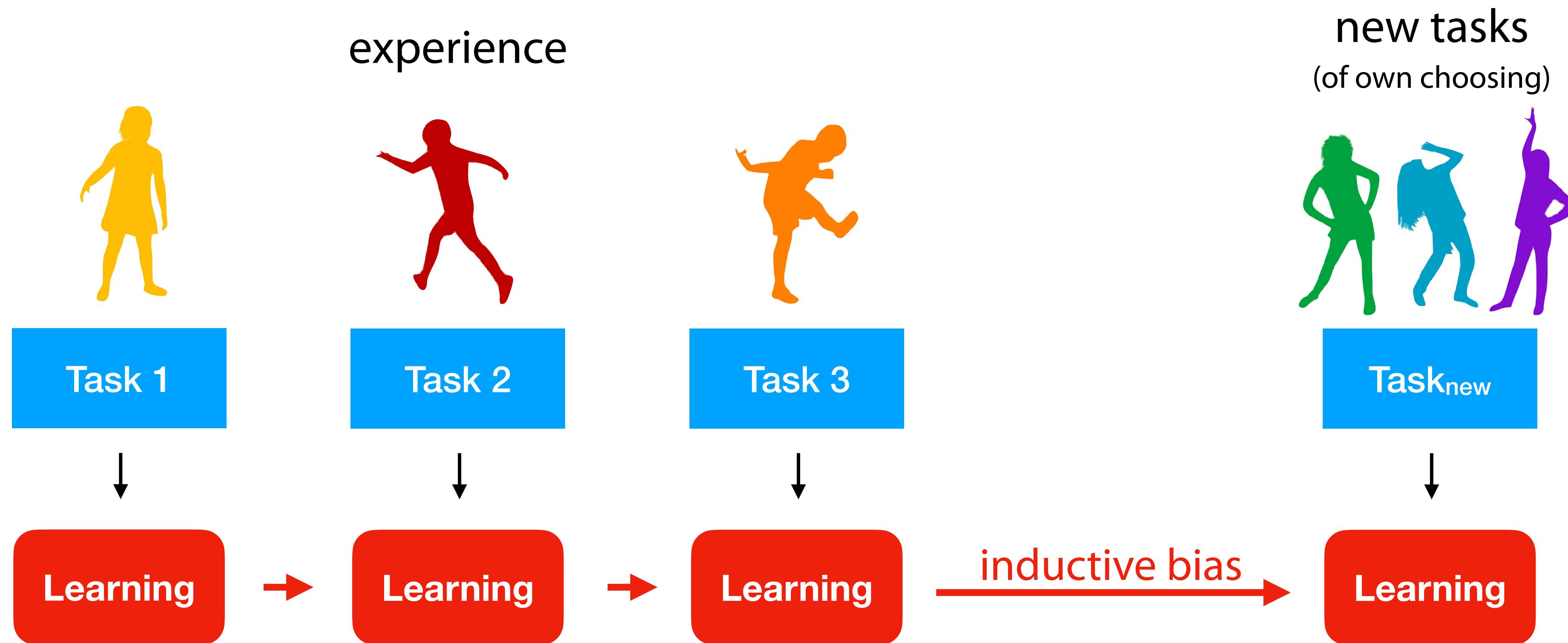


**year 4**



# Human-like Learning\*\*\*

humans learn across tasks: less trial-and-error, less data, less compute  
 new tasks should be related to experience (doable, fun, interesting?)



key aspects of fast learning: compositionality, causality, learning to learn

# Inductive bias (in language)

which assumptions do we make?

Training

- dax
- zup
- lug
- wif

● ● ● lug blicket wif

● ● ● red wif blicket dax

● ● blue lug kiki wif

● ● red wif kiki dax

Test

● ● ● red dax blicket zup?

● ● ● red ● wif blicket dax kiki lug?

Common mistakes



one-to-one bias:  
assume that every word is one color



concatenation bias:  
assume that order is always left-to-right

# What if there is *no* training data?

we can still solve problems by making assumptions

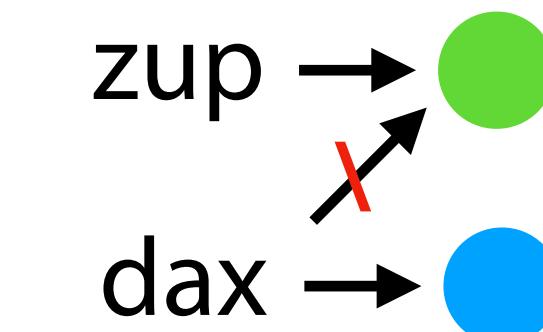
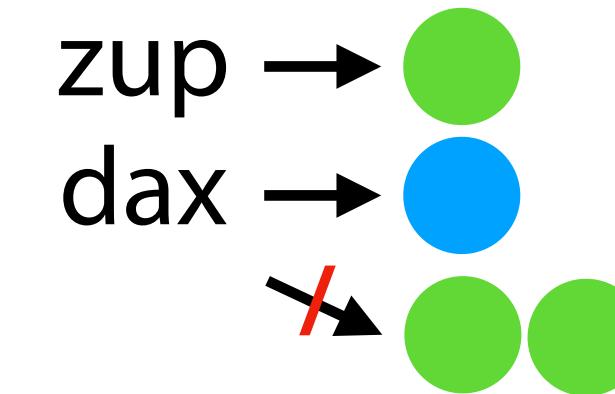
Item pool



Test

- zup?
- zup zup?
- dax zup?
- zup tufa?
- zup wif zup?
- zup wif blicket?
- blicket wif zup?

Commonly used assumptions:



**one-to-one bias:**

assume that every word is one color  
(and not a function or something else)

**concatenation bias:**

assume that order is always left-to-right

**mutual exclusivity:**

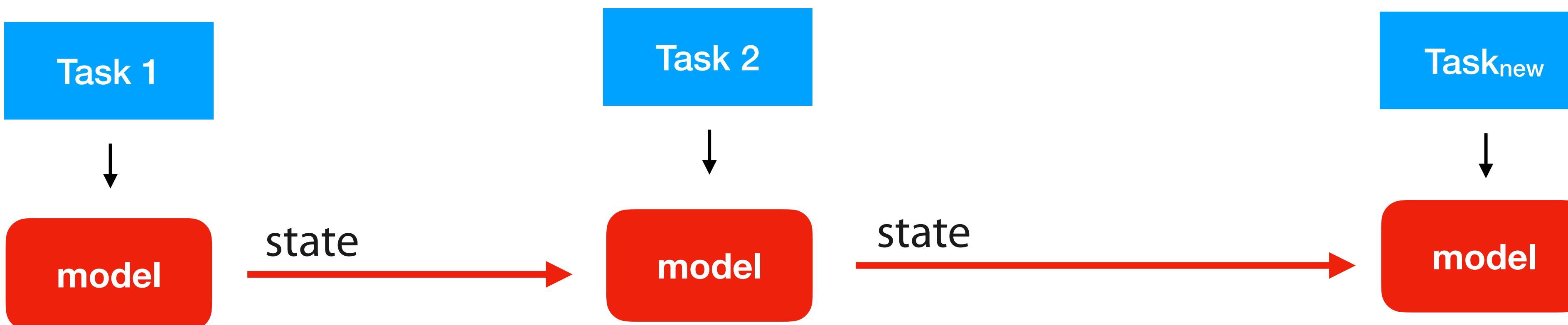
if object has a name, it doesn't need another name

Humans assume that words have consistent meanings and follow input/output constraints

**These assumptions (*inductive biases*) are necessary for learning quickly**

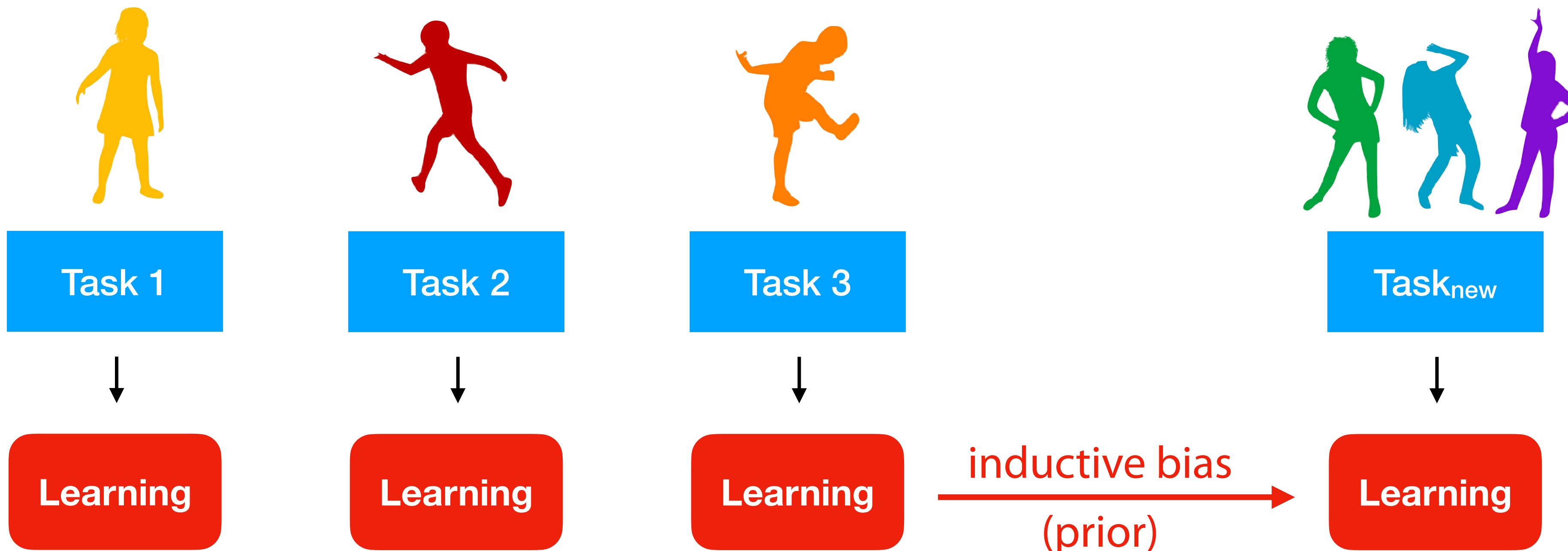
# Meta-learning inductive biases

Capture **useful assumptions** from the data - that can often not be easily expressed



# Meta-learning goal

learn minimal inductive biases from prior tasks instead of constructing manual ones  
should still generalize well (otherwise you meta-overfit)

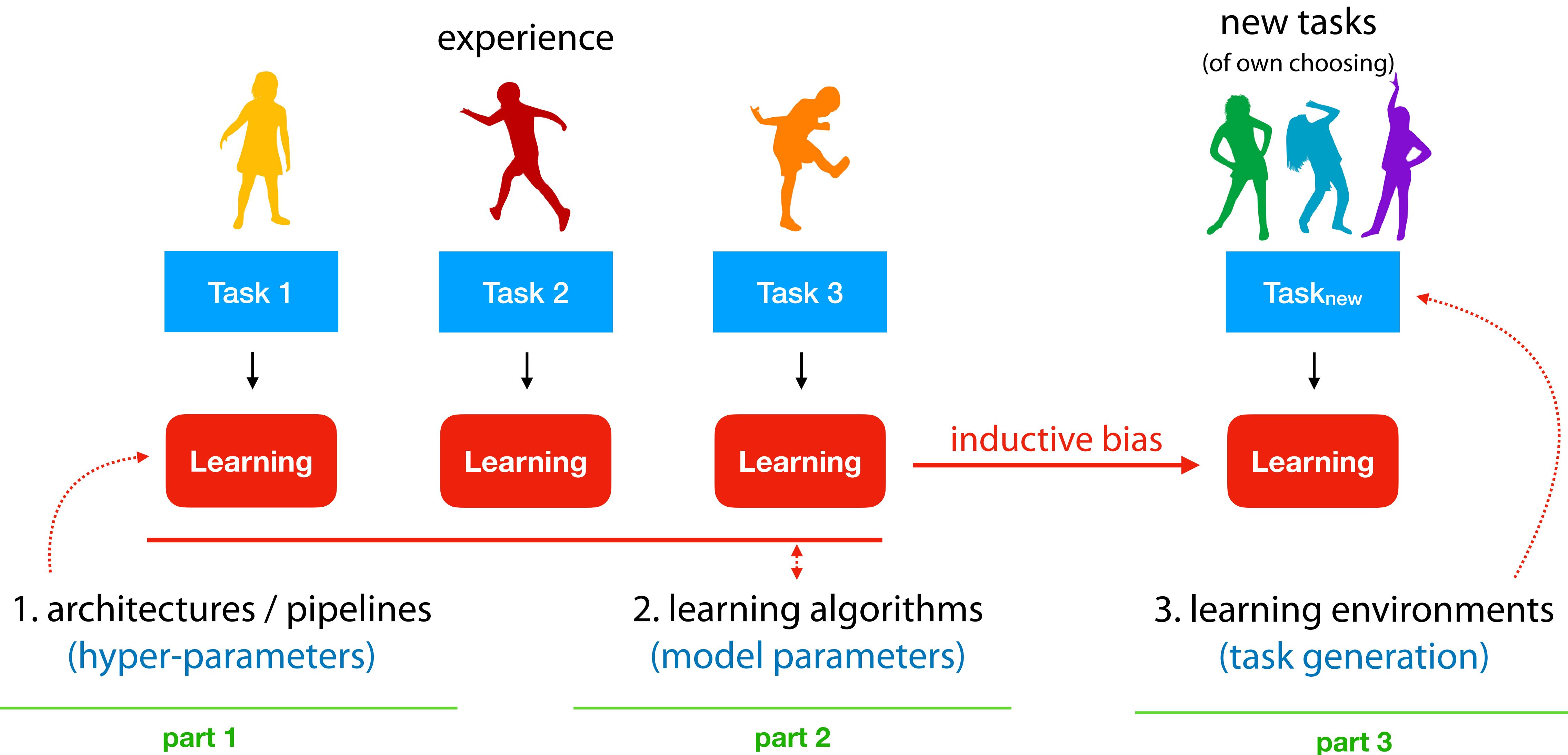


**Inductive bias:** any assumptions added to training data to learn more effectively. E.g:

- Instead of general model architectures, learn better architectures (and hyperparameters)
- Instead of starting from random weights, learn good initial weights
- Instead of standard loss/reward function, learn a better loss/reward function

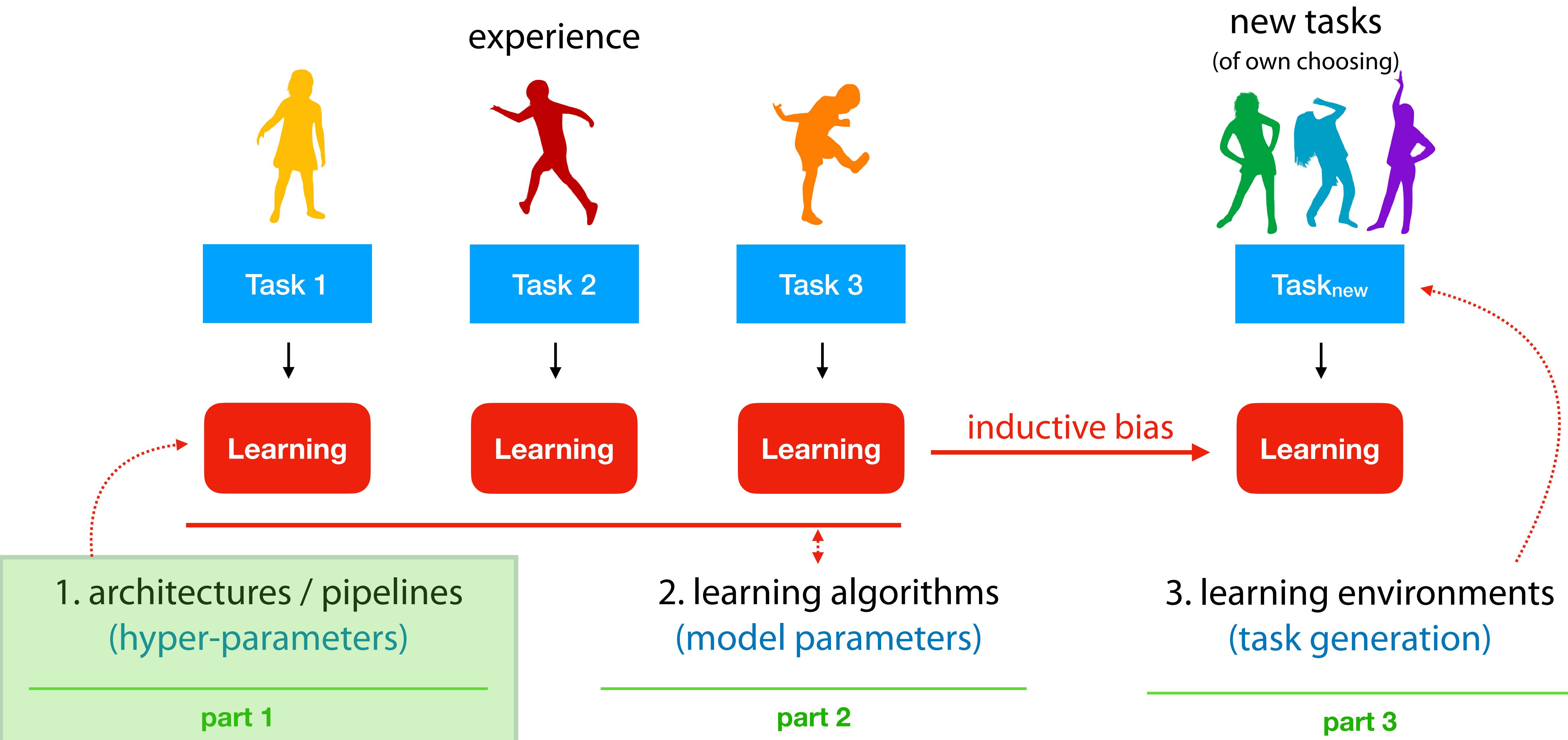
# What can we learn to learn?

*3 pillars*



# What can we learn to learn?

*3 pillars*



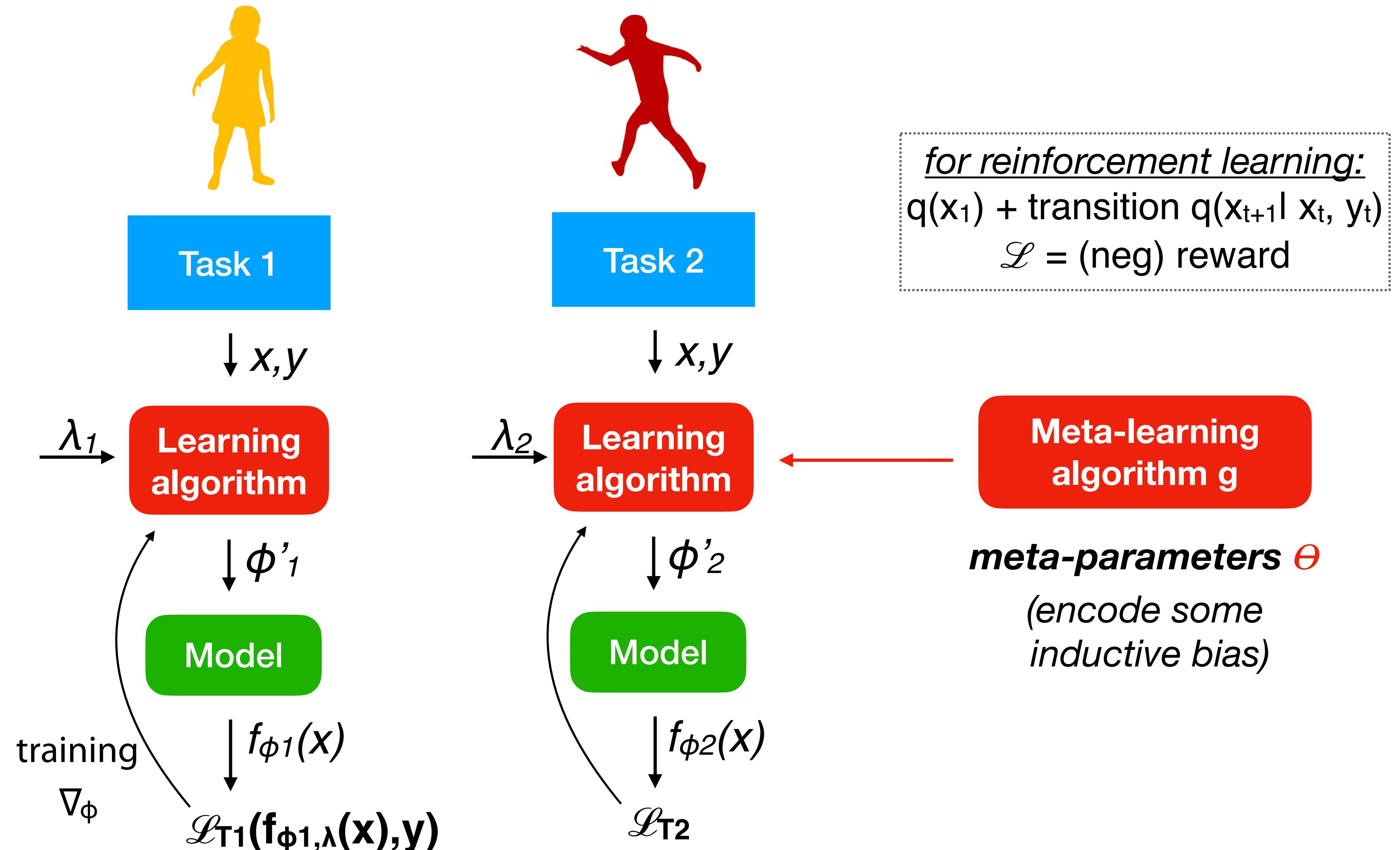
# Terminology

**Task:** distribution of samples  $q(x)$   
outputs  $y$ , loss  $\mathcal{L}(x,y)$

**Learner:** *model parameters*  $\phi$ ,  
*hyper-parameters*  $\lambda$   
optimizer

$$\begin{aligned}\Phi^* &= \operatorname{argmax}_\Phi \log p(\Phi \mid T) \\ &= \operatorname{argmin}_\Phi \mathcal{L}(f_{\Phi,\lambda}(x), y)\end{aligned}$$

**Model:**  $f_\phi(x) = y'$

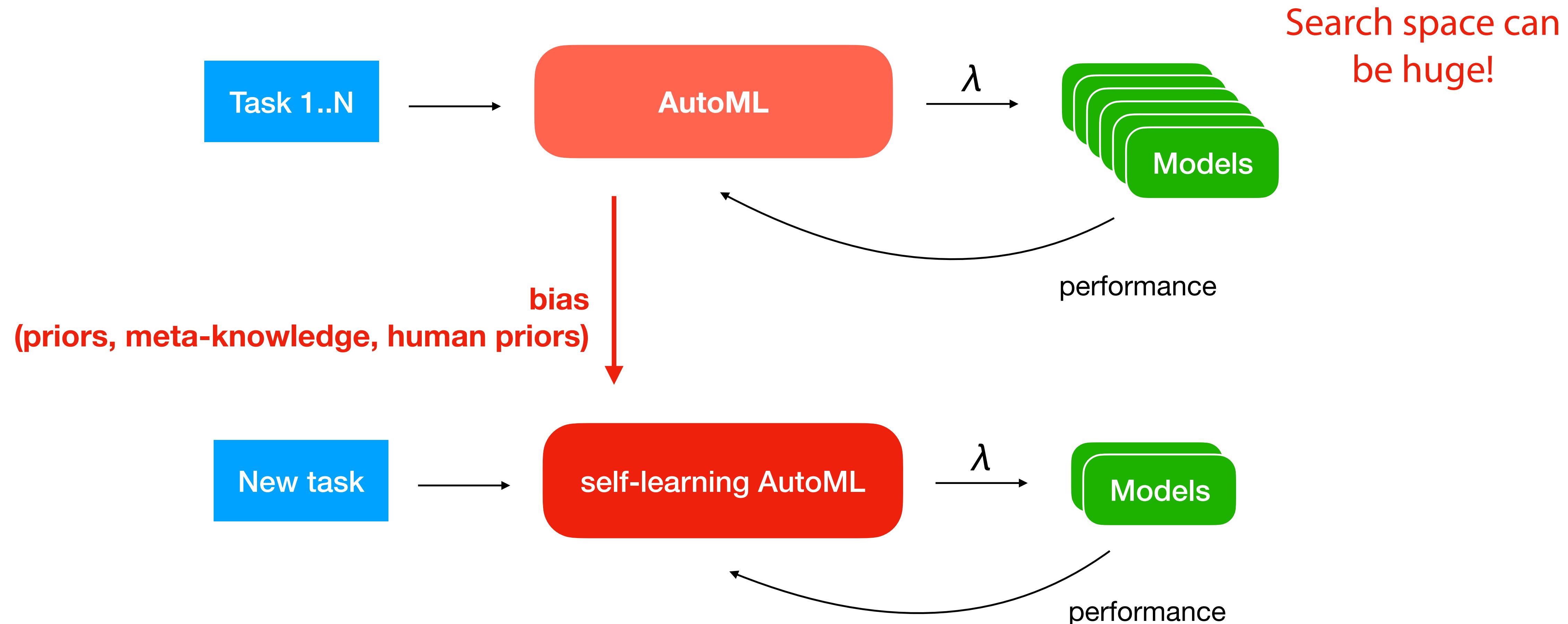


# Learning hyperparameters

Closely related to Automated Machine Learning (AutoML)

But: **meta-learn** how to design architectures/pipelines and tune hyperparameters

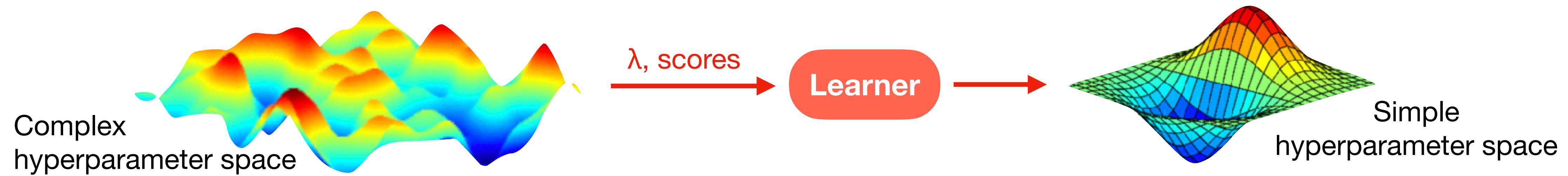
**Human data scientists also learn from experience**



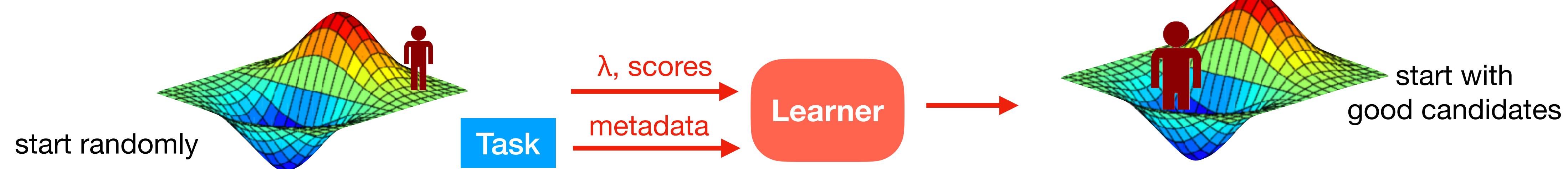
# Meta-learning for AutoML: how?

hyperparameters = architecture + hyperparameters

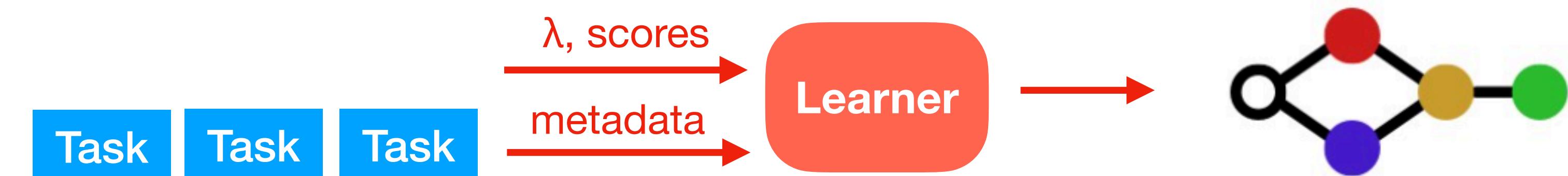
## Learning hyperparameter priors



## Warm starting (what works on similar tasks?)



## Meta-models (learn how to build models/components)



**Observation:**  
**current AutoML strongly depends on learned priors**



autosklearn [Feurer et al. 2015](#)

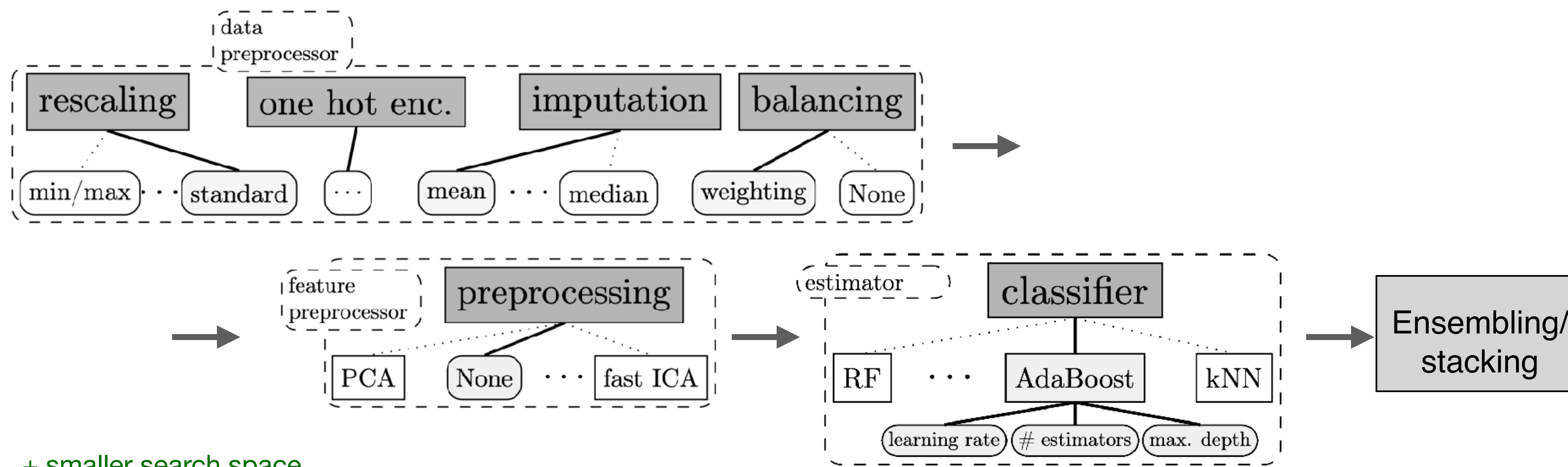
autoWEKA [Thornton et al. 2013](#)

hyperopt-sklearn [Kommer et al. 2014](#)

AutoGluon-Tabular [Erickson et al. 2020](#)

# Manual architecture priors

- Most successful pipelines have a similar structure
- Can we meta-learn a prior over successful structures?



+ smaller search space

- you can't learn entirely new architectures

Figure source: [Feurer et al. 2015](#)

# Manual architecture priors

Successful deep networks often have repeated motifs (cells)

e.g. Inception v4:

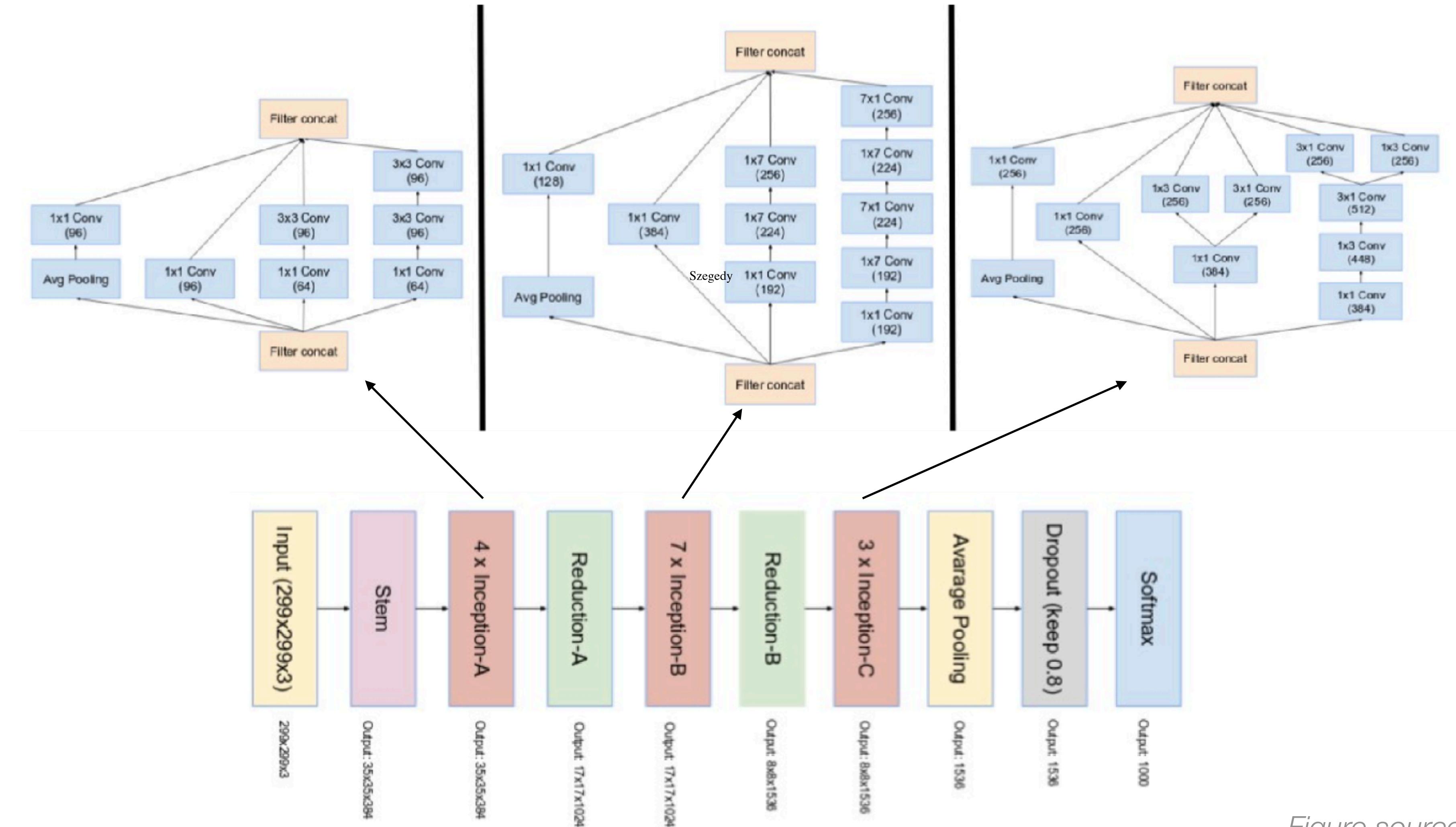


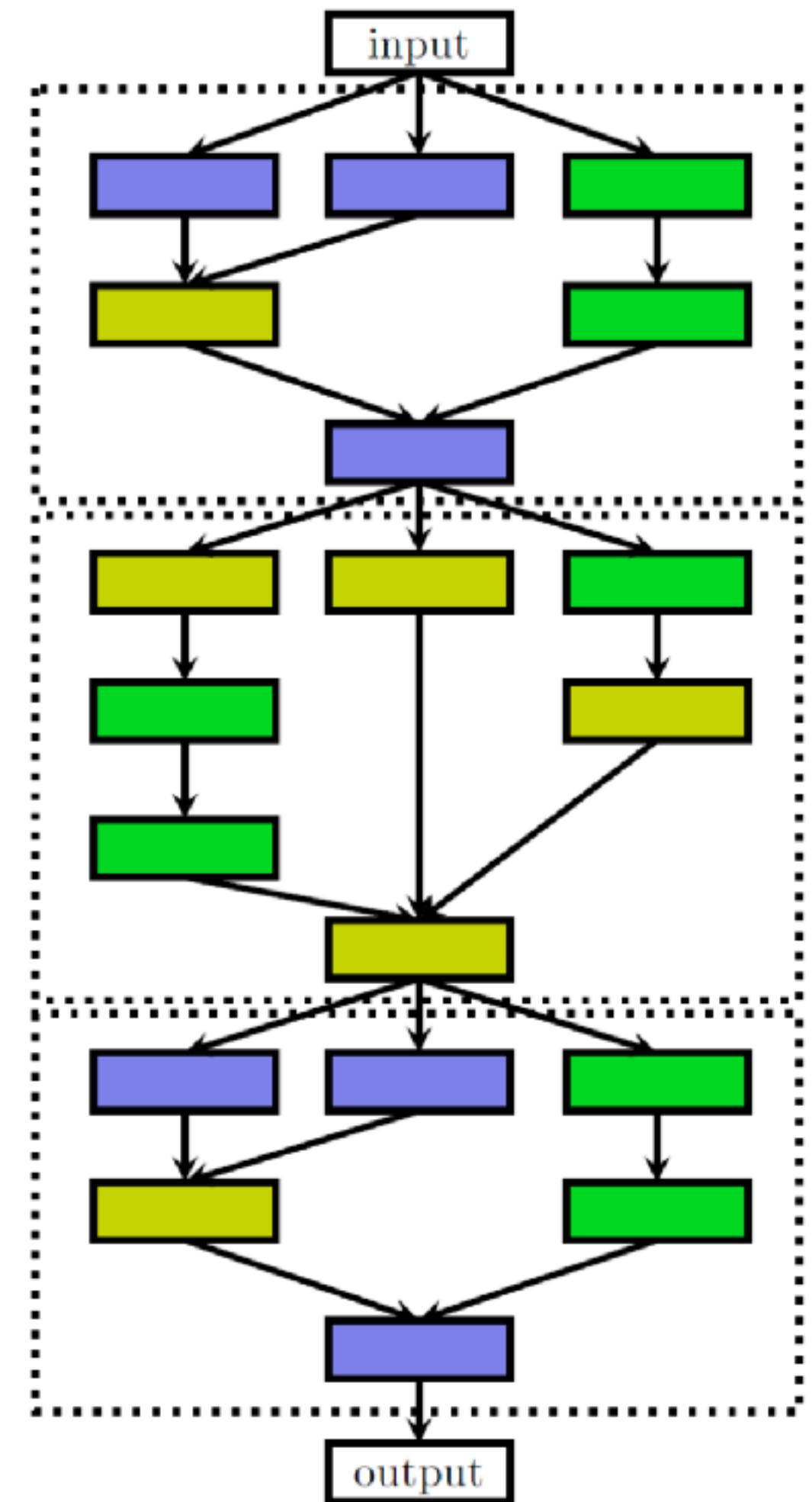
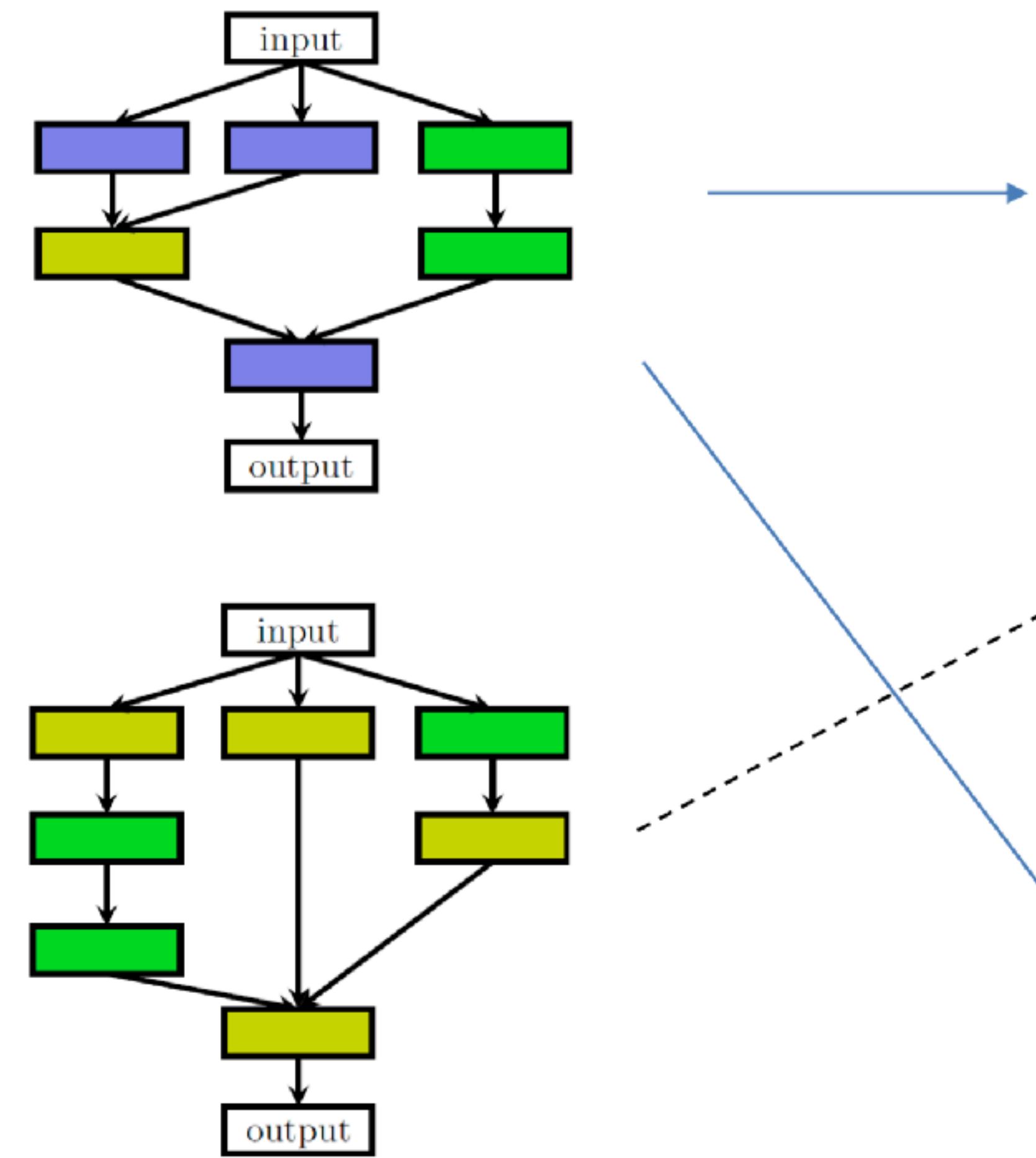
Figure source: Szegedy et al 2016

# Cell search space prior

**Compositionality:** learn hierarchical building blocks to simplify the task

## Cell search space

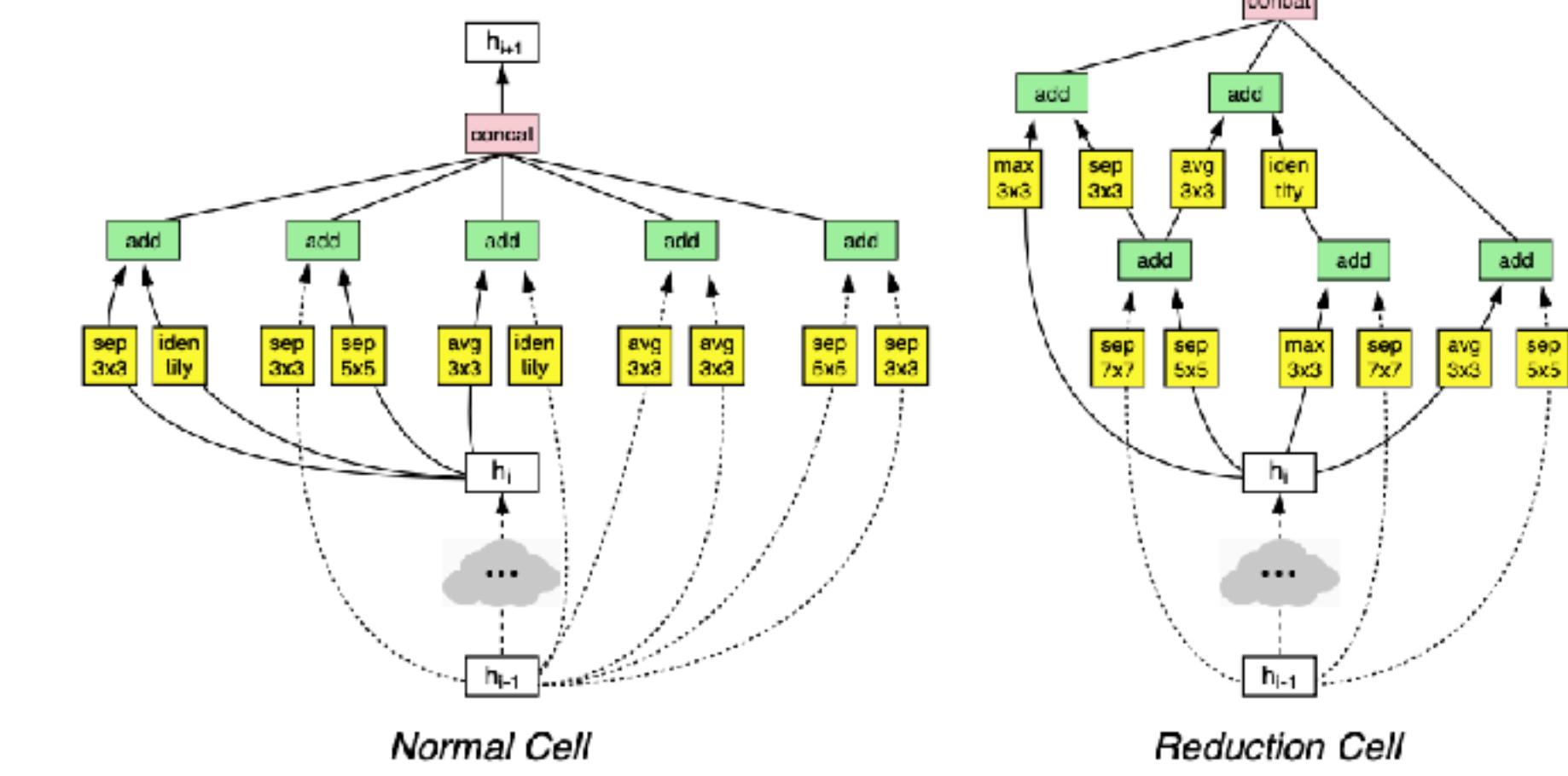
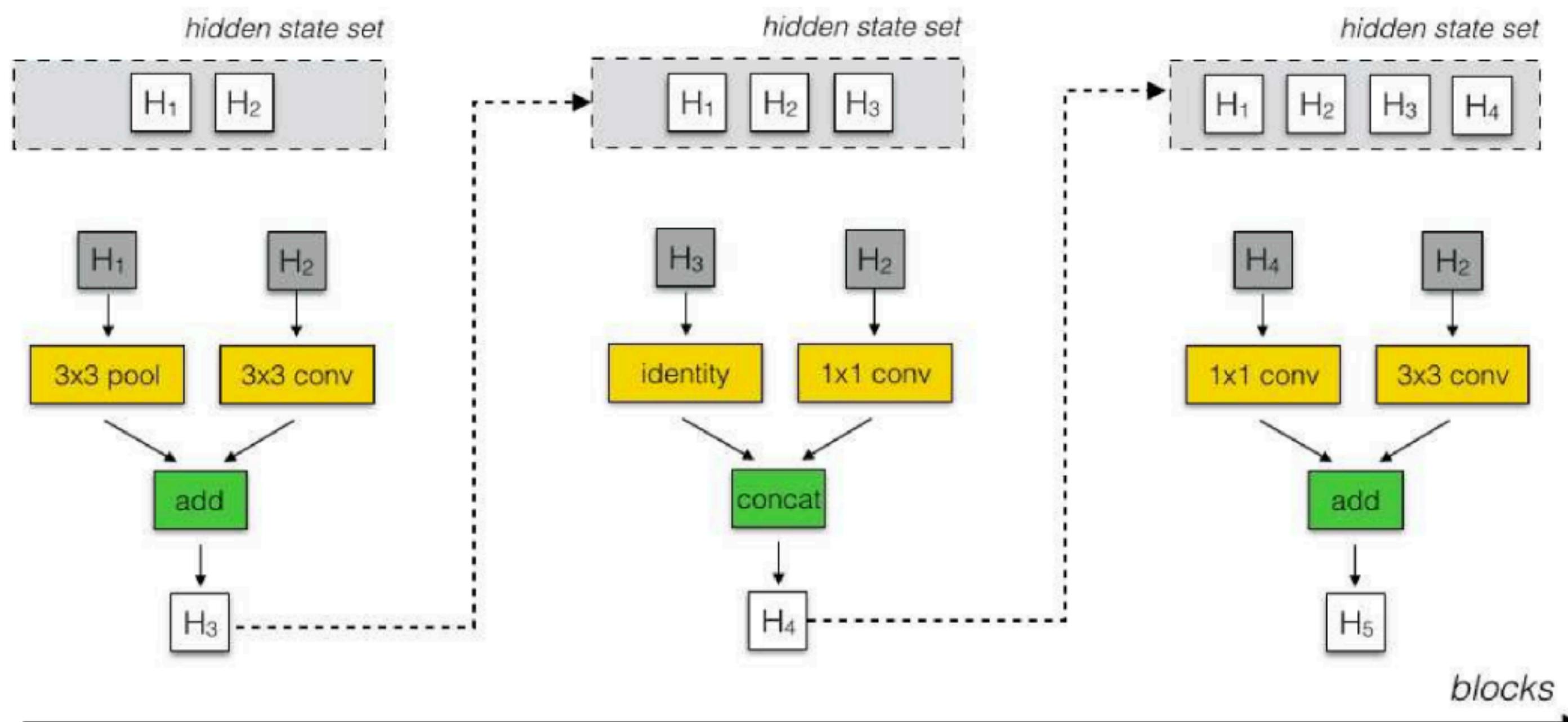
- learn parameterized building blocks (*cells*)
  - stack cells together in macro-architecture
- + smaller search space  
+ cells can be learned on a small dataset & transferred to a larger dataset  
- strong domain priors, doesn't generalize well



Can we meta-learn hierarchies / components that generalize better?

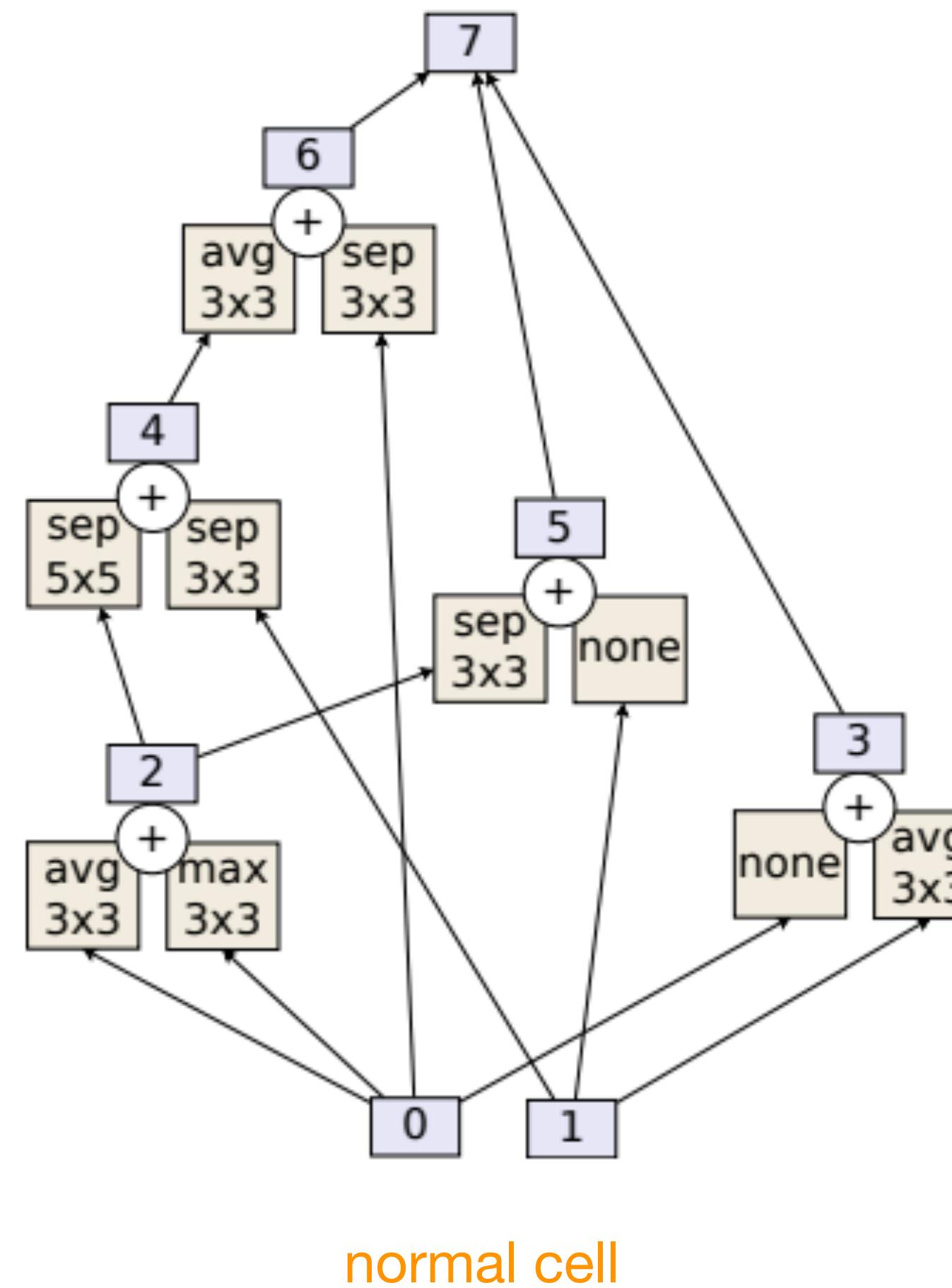
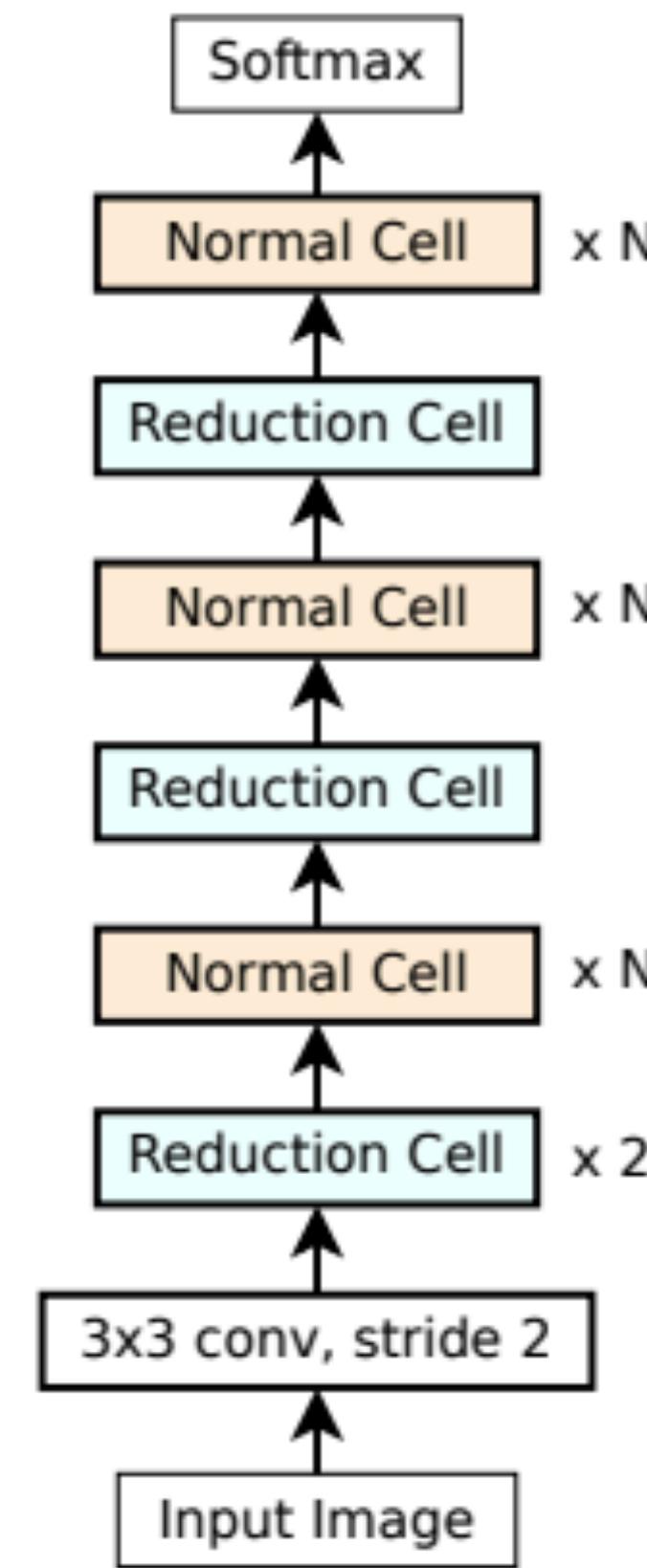
# Cell search space prior

- Cell construction with **RL-based search** (SOTA ImageNet, **450 GPUs, 3-4 days**):
  - Select existing layers (hidden states, e.g. cell input)  $H_i$  to build on
  - Add operation (e.g. 3x3conv) on  $H_i$
  - Combine into new hidden state (e.g. concat, add,...)
  - Iterate over  $B$  blocks

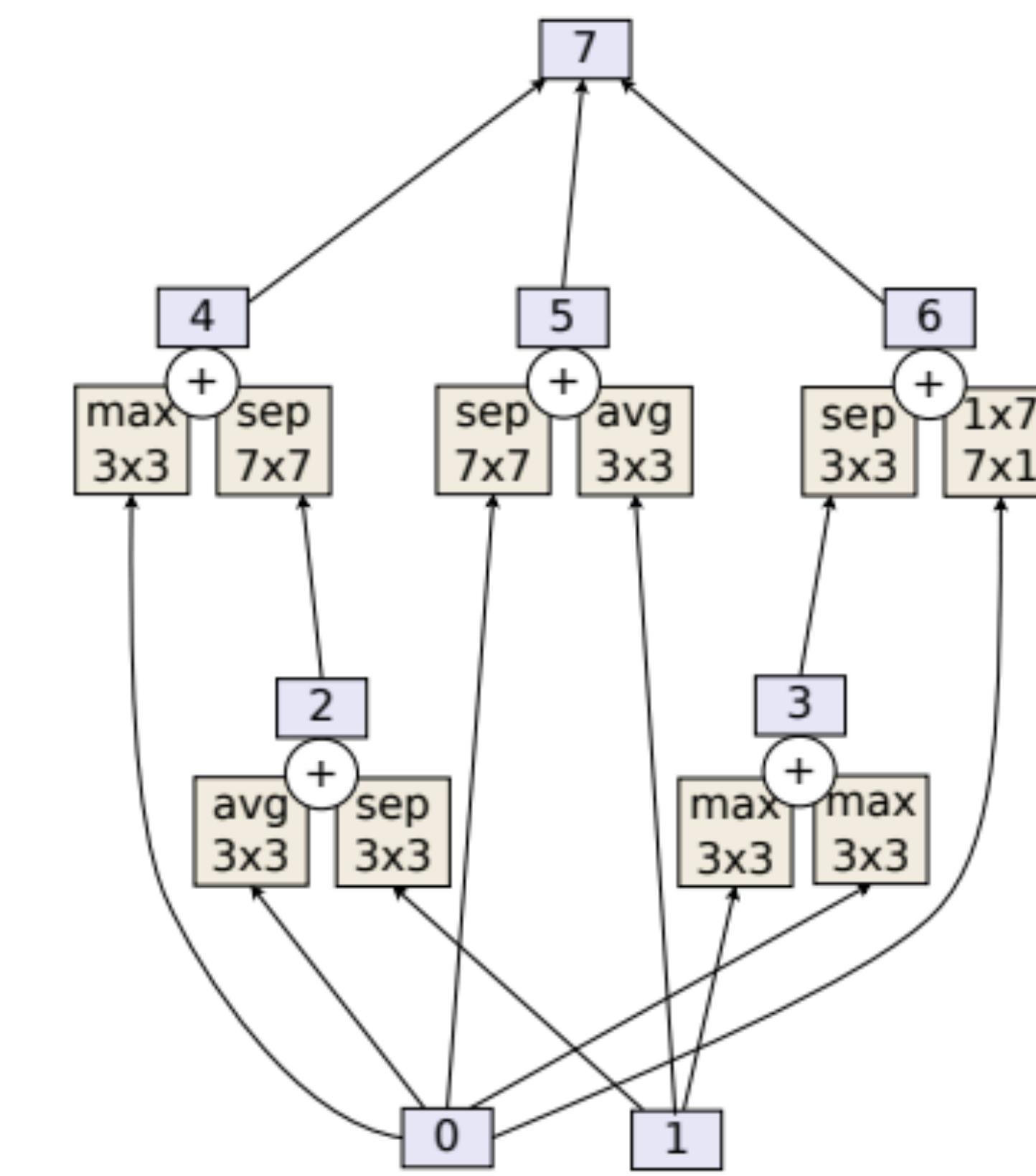


# Cell search space prior

- Cell construction with **neuro-evolution** (SOTA ImageNet)



normal cell

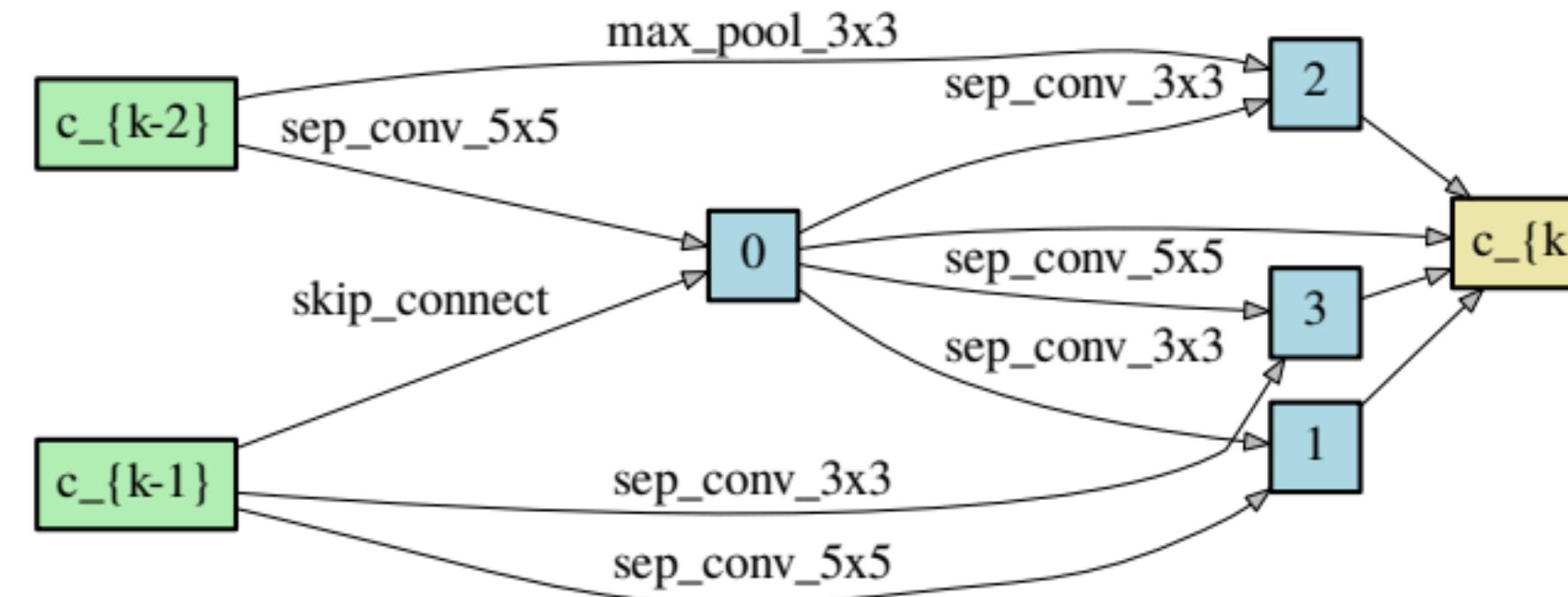


reduction cell

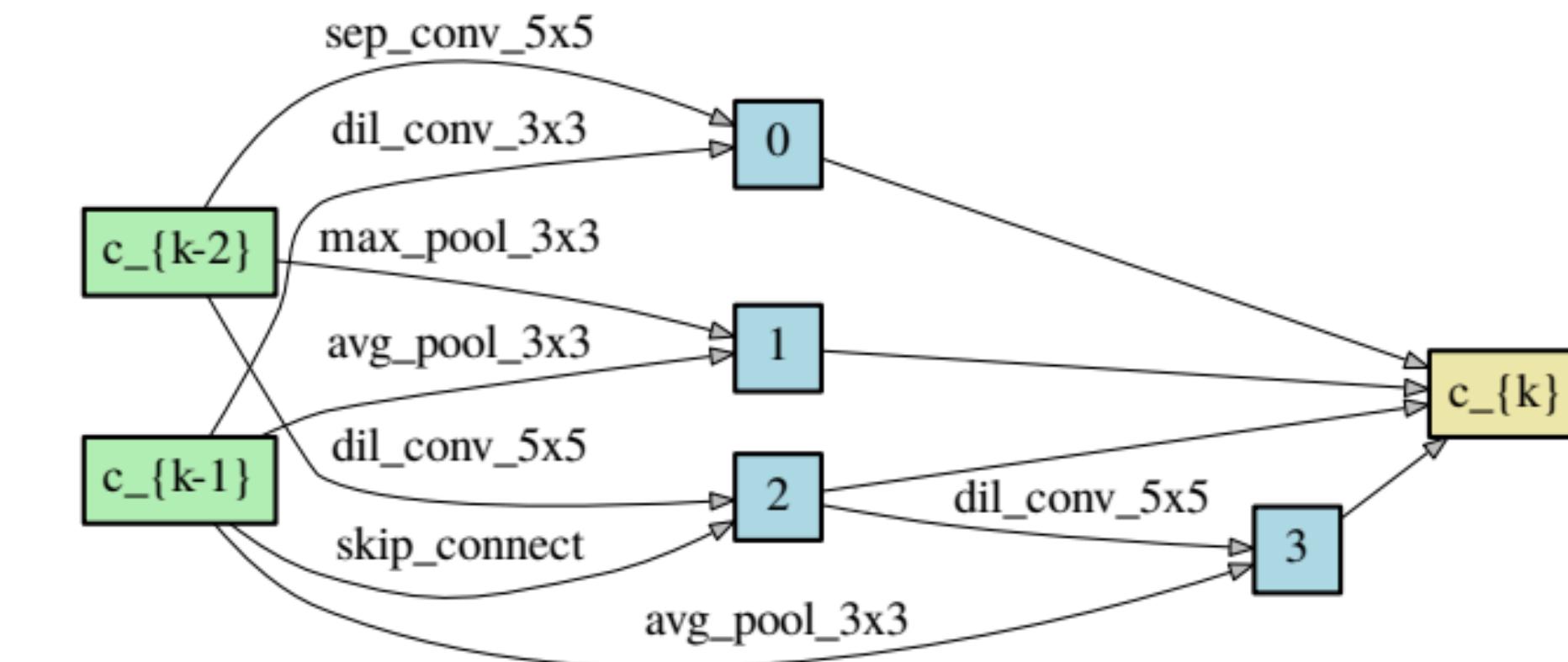
# Cell search space prior

- Cell construction with **multi-fidelity random search!**

*If you constrain the search space enough, you can get SOTA results with random search!*



(a) Normal Cell



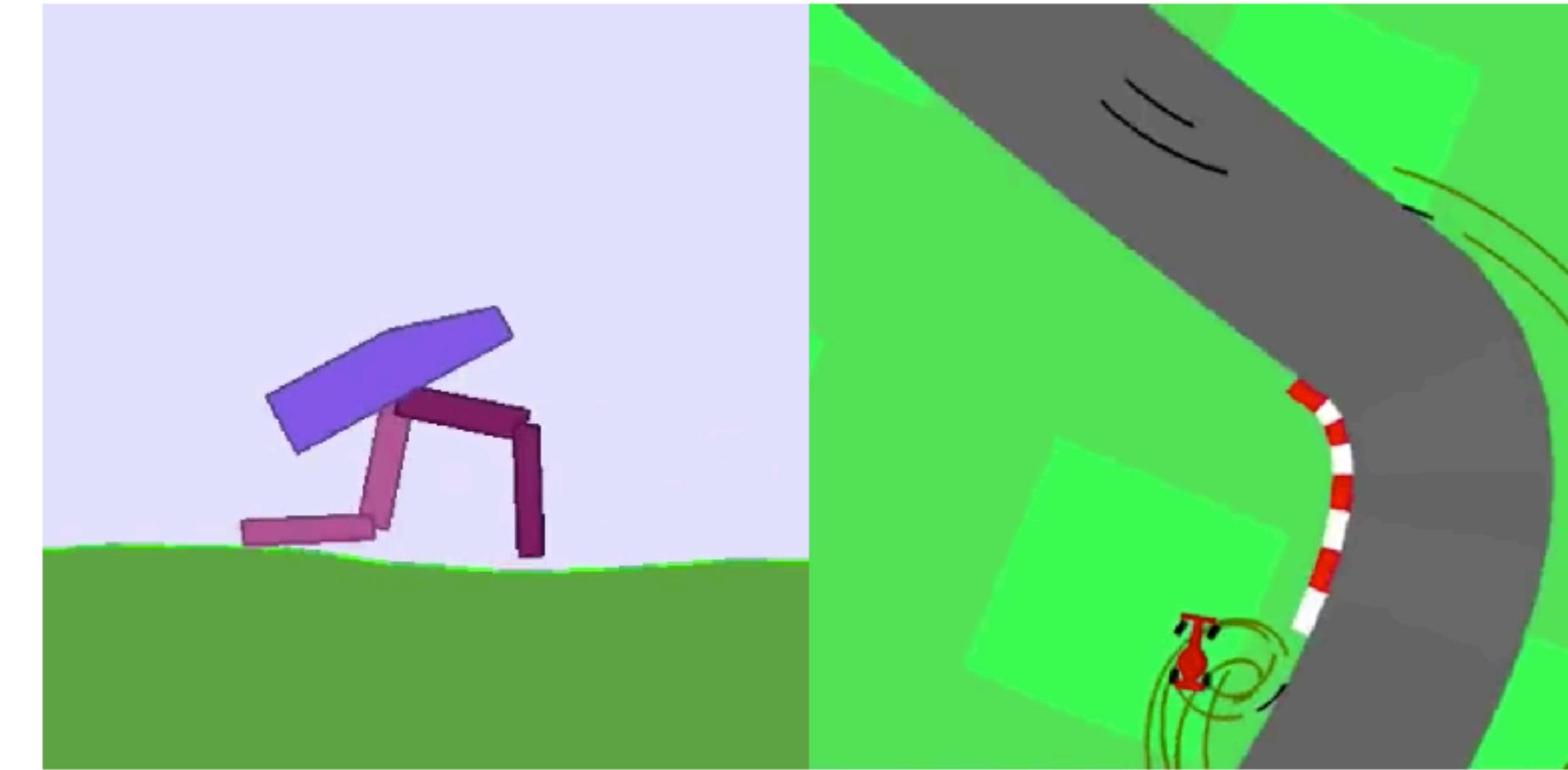
(b) Reduction Cell

**Convolutional Cells on CIFAR-10 Benchmark:** Best architecture found by random search with weight-sharing.

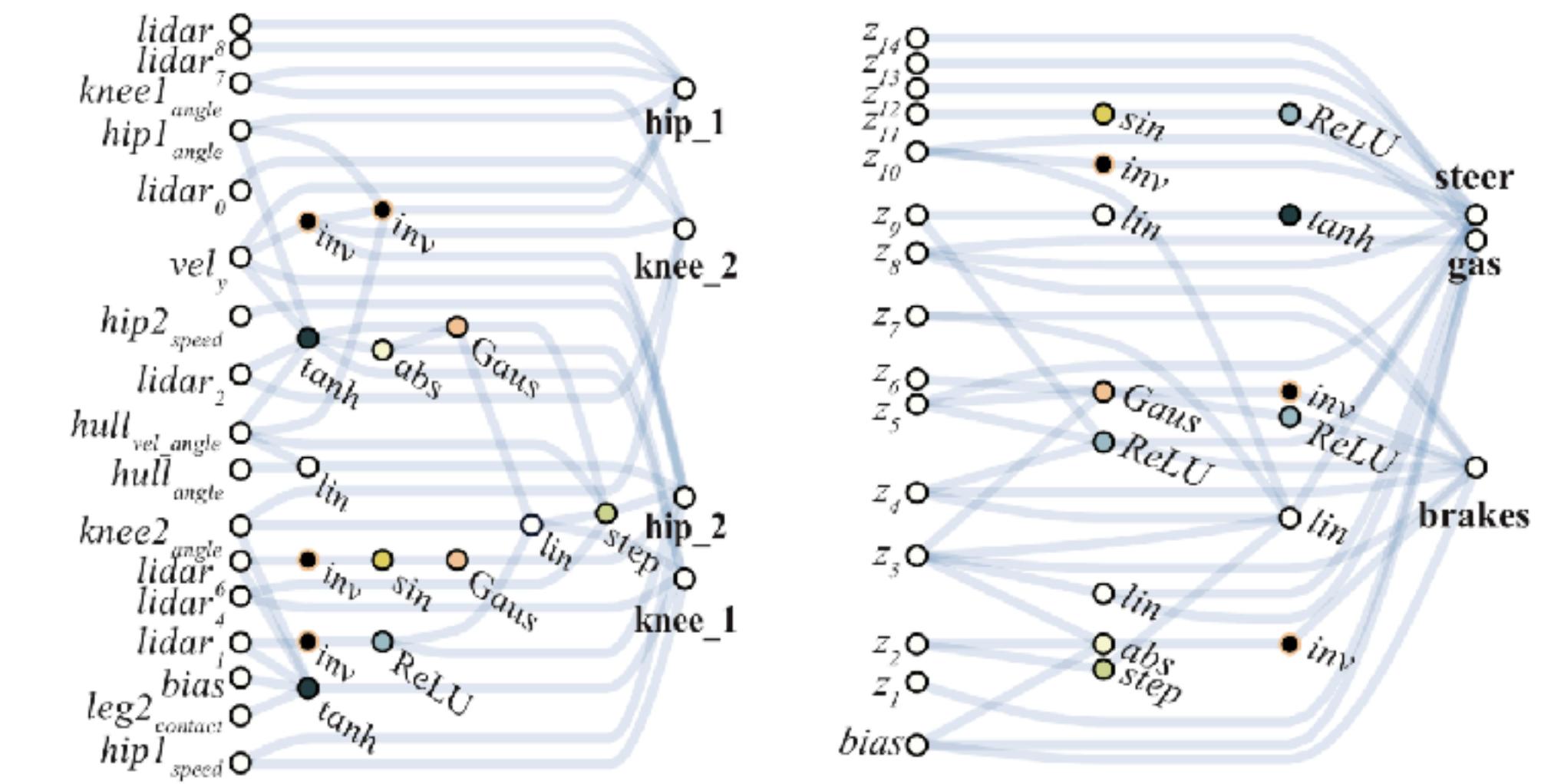
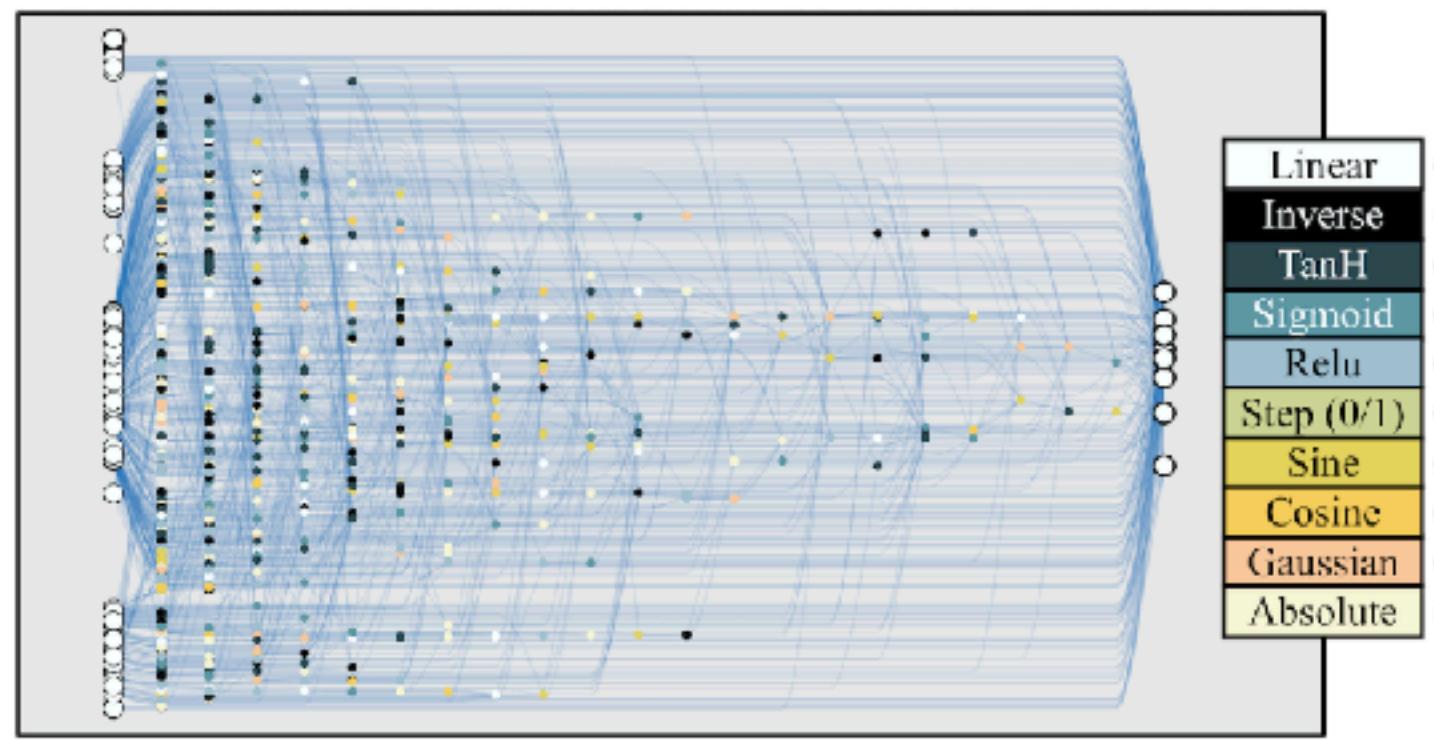
# Manual priors: Weight sharing

## Weight-agnostic neural networks

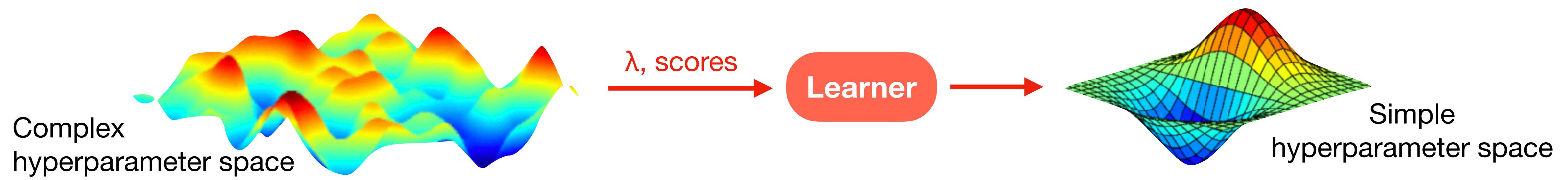
- ALL weights are shared
- Only evolve the architecture?
  - Minimal description length
  - Baldwin effect?



MNIST digit



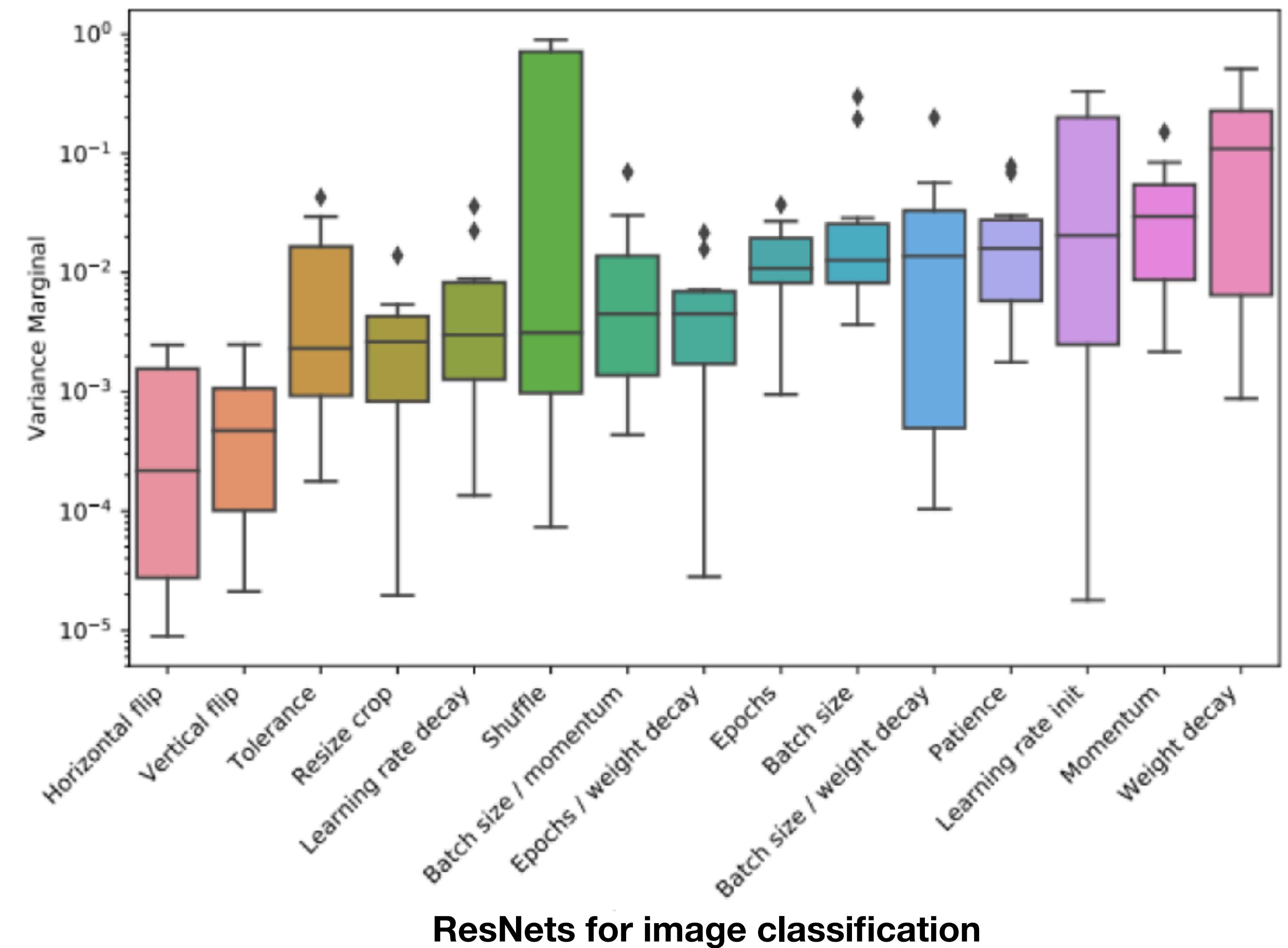
# Learning hyperparameter priors



# Learn hyperparameter importance

- **Functional ANOVA** <sup>1</sup>

- Select hyperparameters that cause variance in the evaluations.
- Useful to speed up black-box optimization techniques



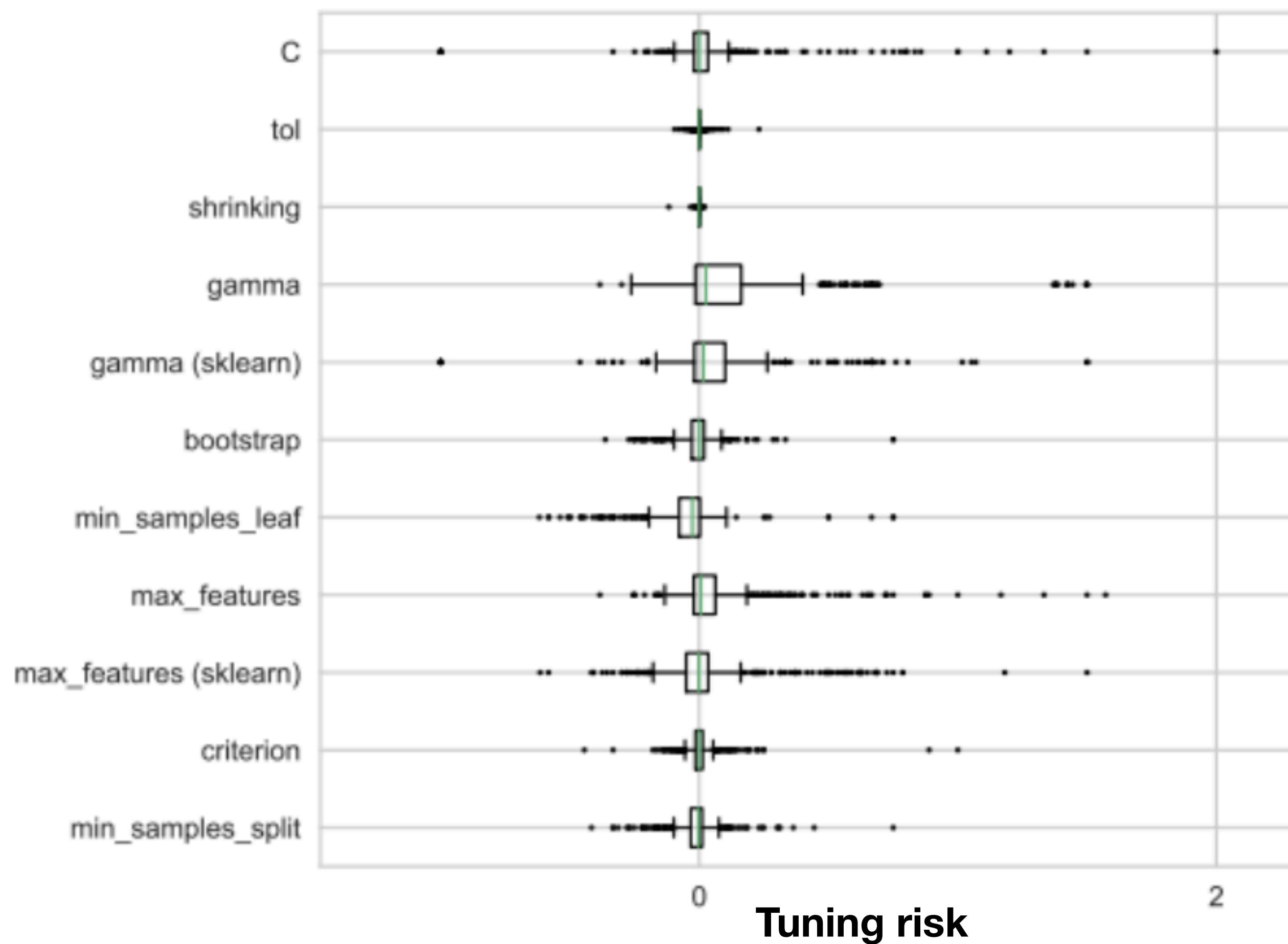
# Learn defaults + hyperparameter importance

- **Tunability**<sup>1,2,3</sup>

**Learn** good defaults, measure importance as **improvement** via tuning

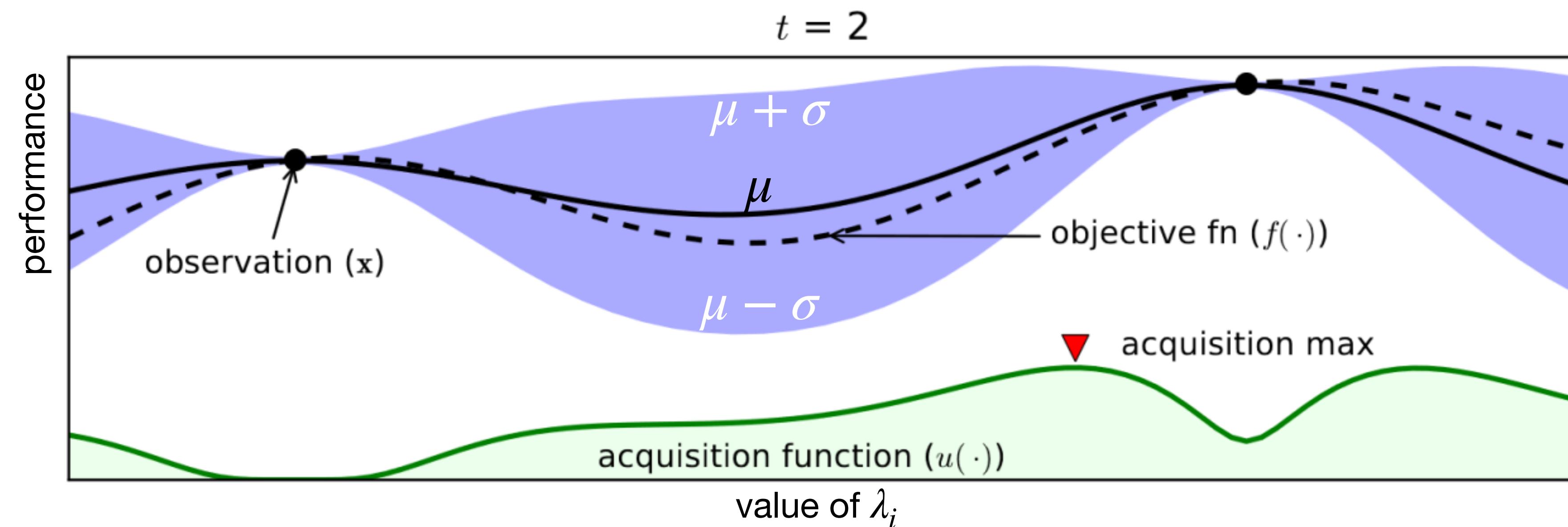
function
max_features
$m = 0.16^*p$
$m = p^{0.74}$
$m = 1.15^{\sqrt{p}}$
$m = \sqrt{p}$
gamma
$m = 0.00574^*p$
$m = 1/p$
$m = 0.006$

Learned defaults



# Bayesian Optimization (interlude)

- Start with a few (random) hyperparameter configurations
- Fit a **surrogate model** to predict other configurations
- Probabilistic regression: mean  $\mu$  and standard deviation  $\sigma$  (blue band)
- Use an **acquisition function** to trade off exploration and exploitation, e.g. Expected Improvement (EI)
- Sample for the best configuration under that function



# Bayesian Optimization

- Repeat until some stopping criterion:
  - Fixed budget
  - Convergence
  - EI threshold
- Theoretical guarantees  
Srinivas et al. 2010, Freitas et al. 2012, Kawaguchi et al. 2016
- Also works for non-convex, noisy data
- Used in AlphaGo

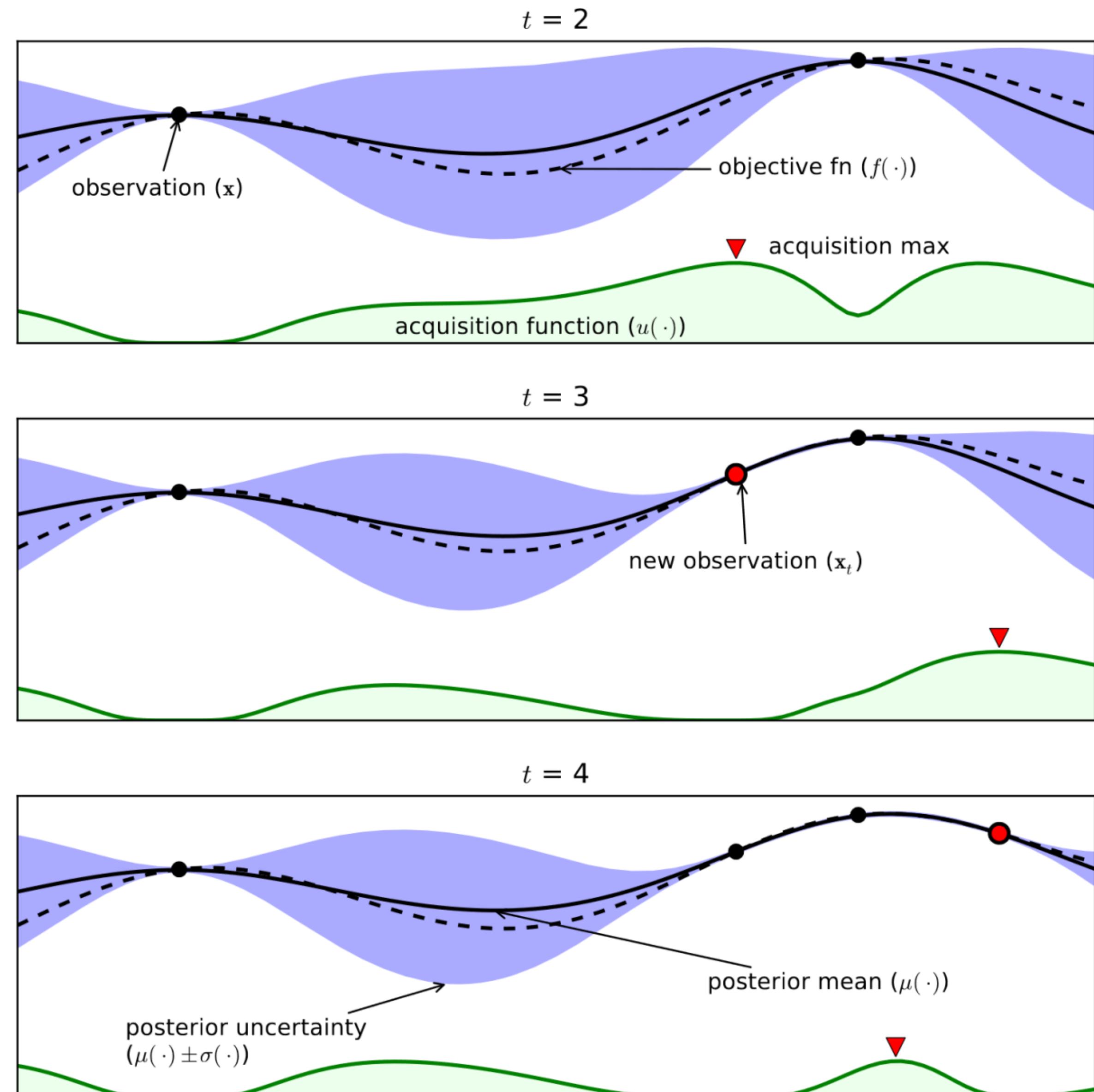
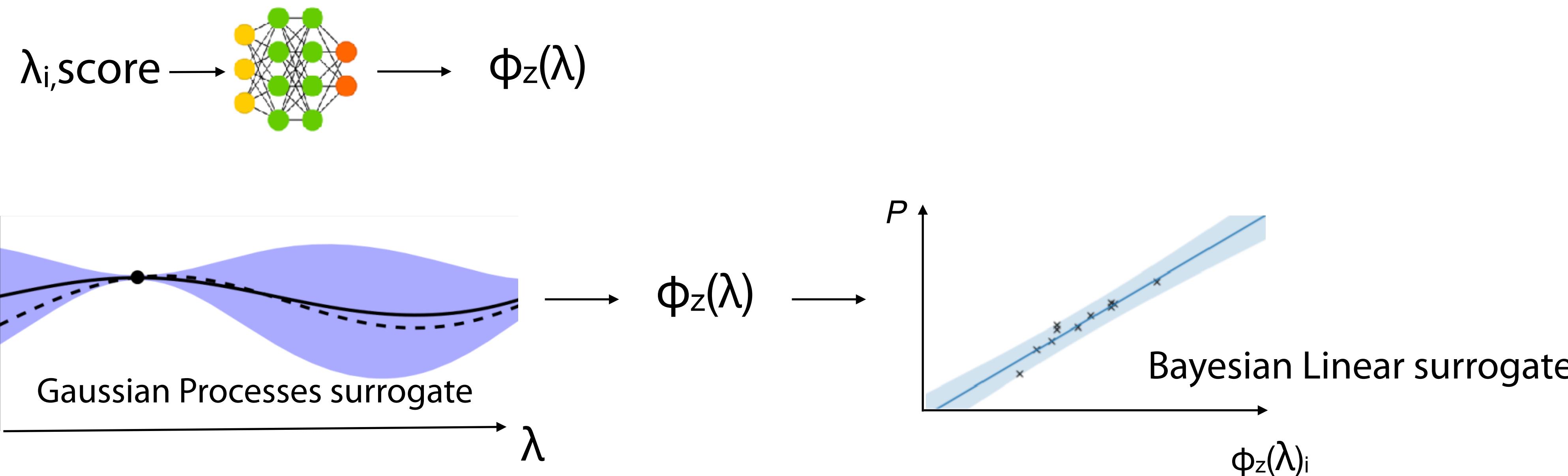


Figure source: Shahriari 2016

# Learn basis expansions for hyperparameters

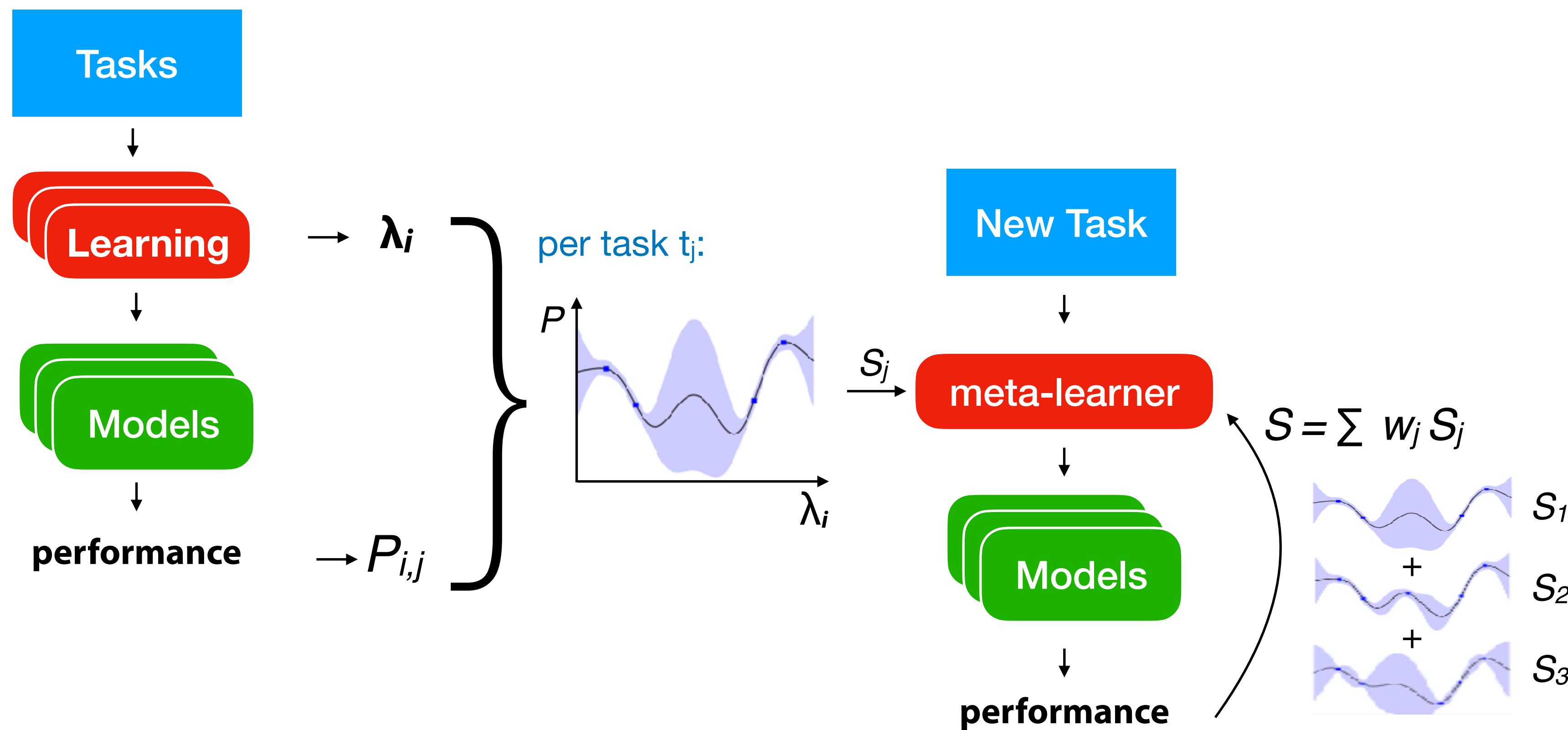
- Hyperparameters can interact in very non-linear ways
- Use a neural net to learn a suitable basis expansion  $\phi_z(\lambda)$  for all tasks
- You can use Bayesian linear models, transfers info on configuration space

Learn basis expansion on lots of data (e.g. OpenML)



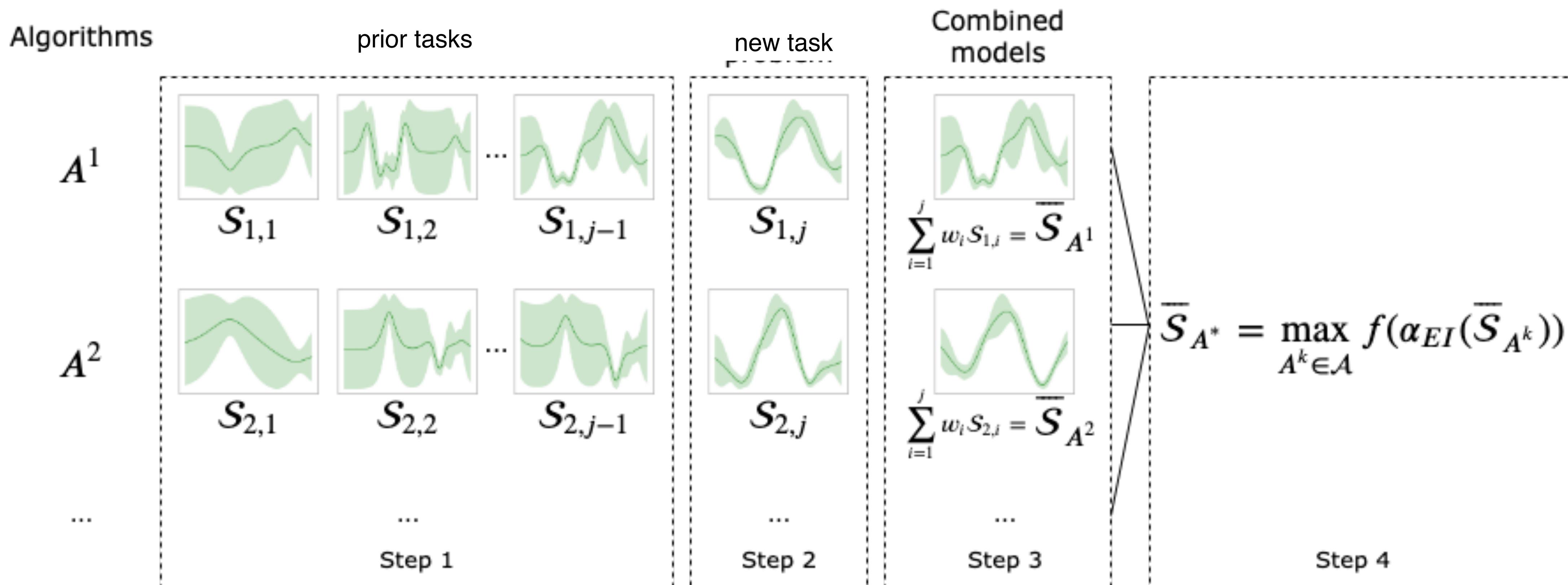
# Surrogate model transfer

- If task  $j$  is similar to the new task, its surrogate model  $S_j$  will likely transfer well
- Sum up all  $S_j$  predictions, weighted by task similarity (as in active testing)<sup>1</sup>
- Build combined Gaussian process, weighted by current performance on new task<sup>2</sup>



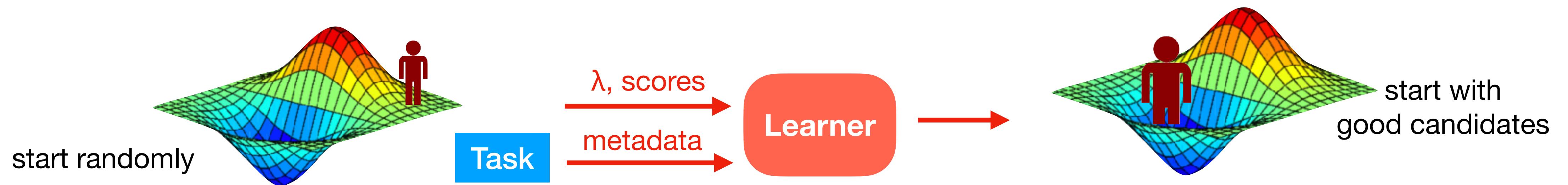
# Surrogate model transfer

- Store surrogate model  $S_{ij}$  for every pair of task i and algorithm j
- Simpler surrogates, better transfer
- Learn weighted ensemble -> significant speed up in optimization



# Warm starting

(what works on similar tasks?)



<sup>1</sup> [Vanschoren 2018](#)

<sup>2</sup> [Achille et al. 2019](#)

<sup>3</sup> [Alvarez-Melis et al. 2020](#)

<sup>4</sup> [Drori et al. 2019](#)

<sup>5</sup> [Jooma et al. 2020](#)

<sup>6</sup> [de Bie et al. 2020](#)

# How to measure task similarity?

- Hand-designed (statistical) meta-features that describe (tabular) datasets <sup>1</sup>
- Task2Vec: task embedding for image data <sup>2</sup>
- Optimal transport: similarity measure based on comparing probability distributions <sup>3</sup>
- Metadata embedding based on textual dataset description <sup>4</sup>
- Dataset2Vec: compares batches of datasets <sup>5</sup>
- Distribution-based invariant deep networks <sup>6</sup>

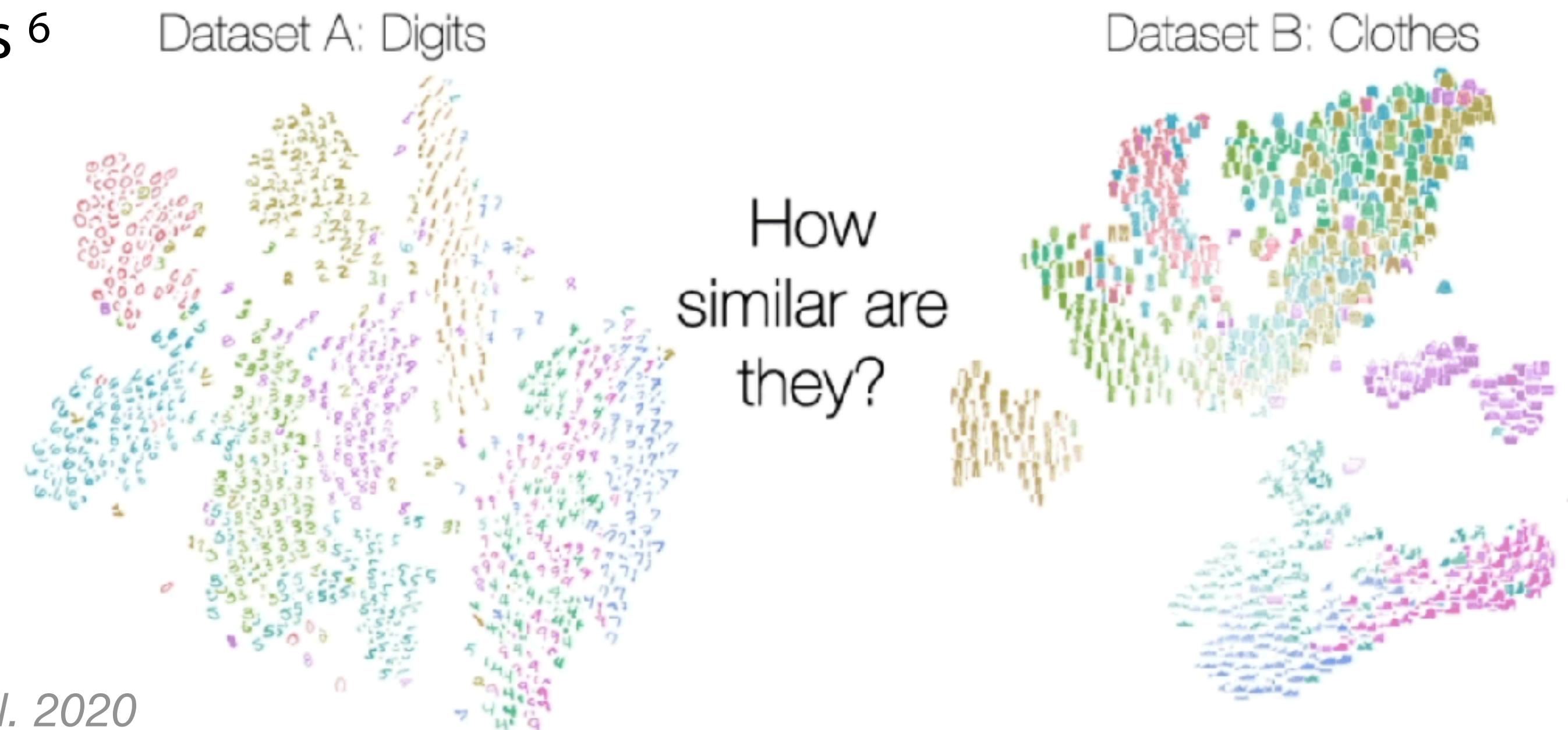
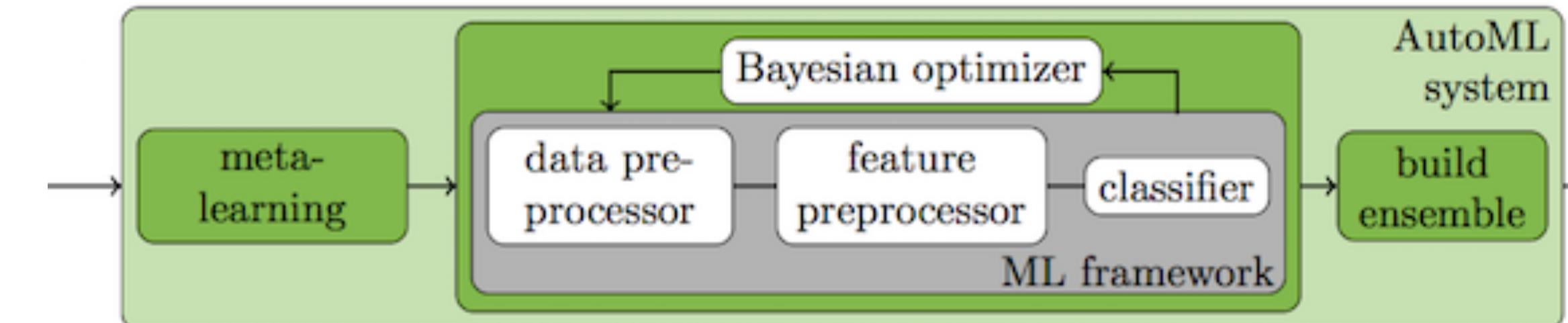


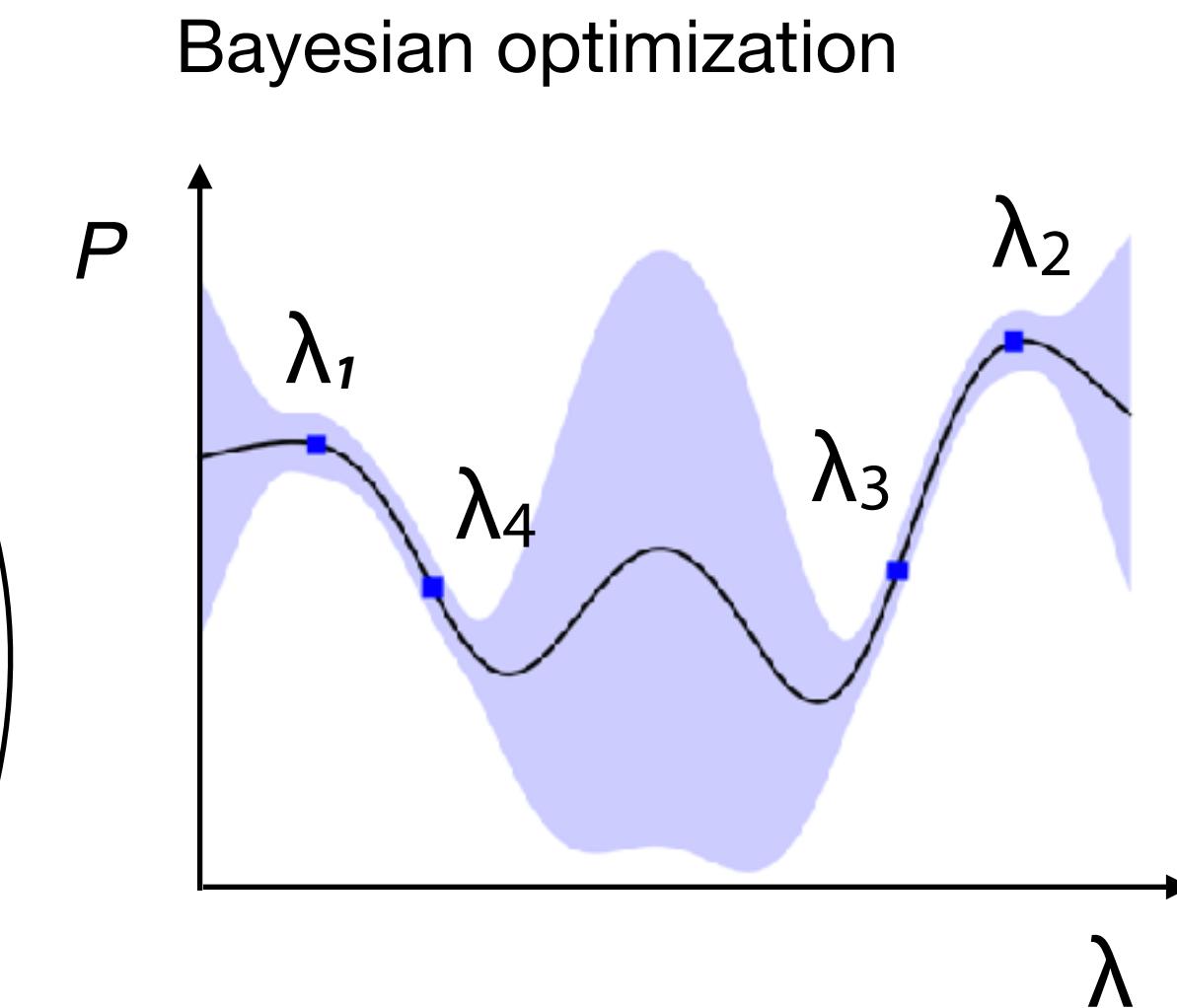
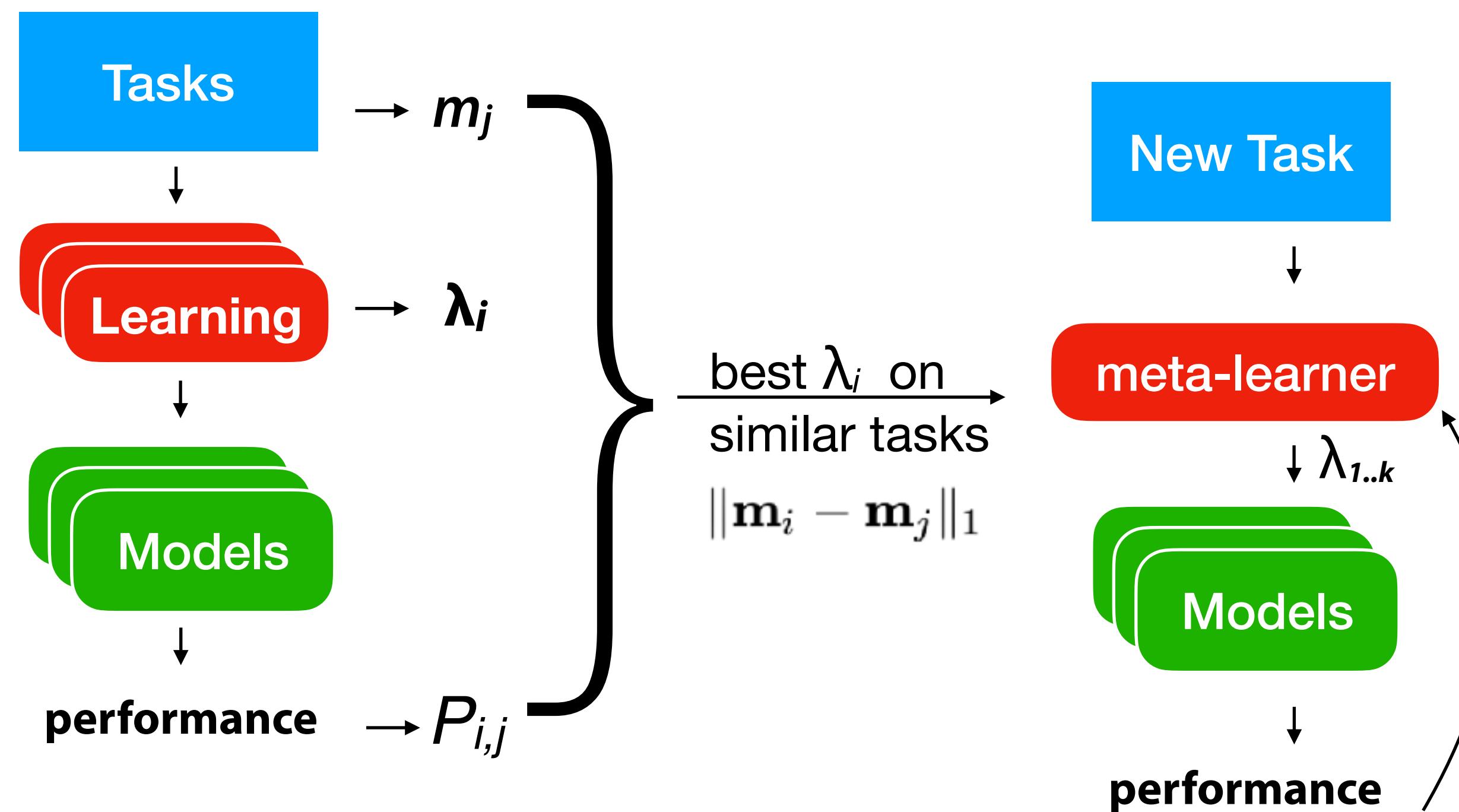
Figure source: Alvarez-Melis et al. 2020

# Warm-starting with kNN

- Find  $k$  most similar tasks, warm-start search with best  $\lambda_i$ 
  - Auto-sklearn: Bayesian optimization (SMAC)
    - Meta-learning yield better models, faster
    - Winner of AutoML Challenges

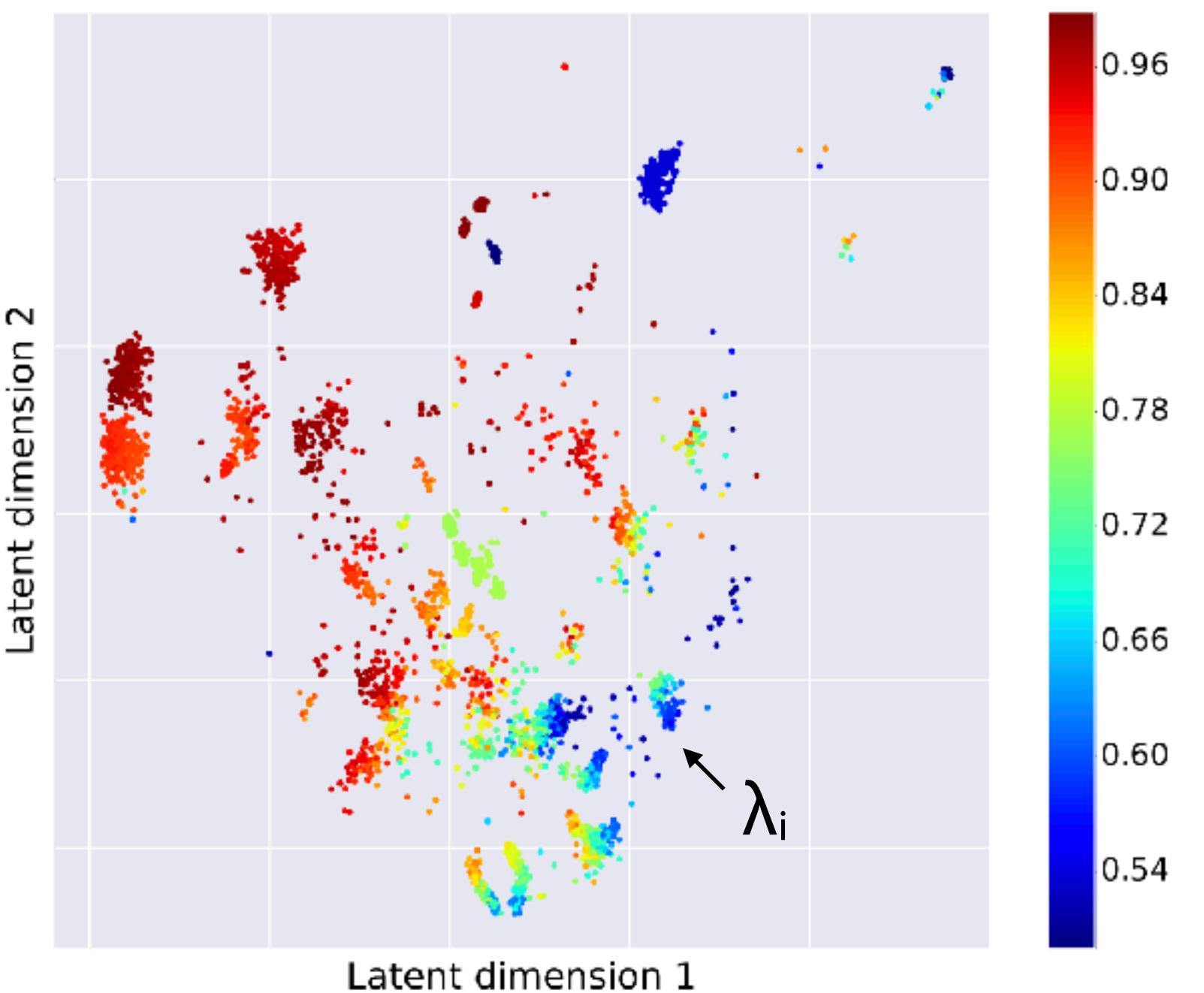
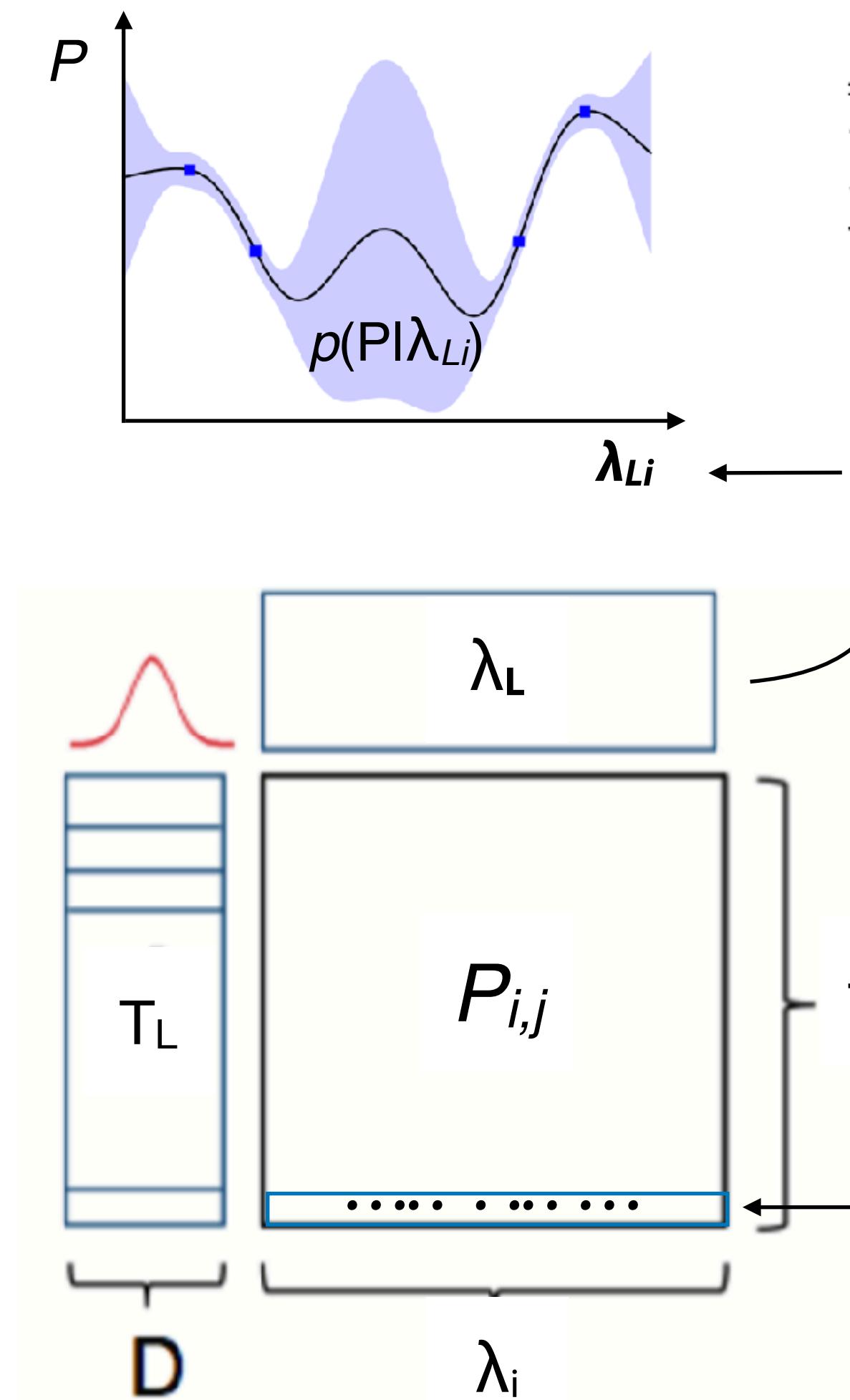


*Figure source: Feurer et al., 2015*



# Probabilistic Matrix Factorization

- Collaborative filtering: configurations  $\lambda_i$  are ‘rated’ by tasks  $t_j$
- Learn latent representation for tasks  $T$  and configurations  $\lambda$
- Use meta-features to warm-start on new task
- Returns probabilistic predictions for Bayesian optimization



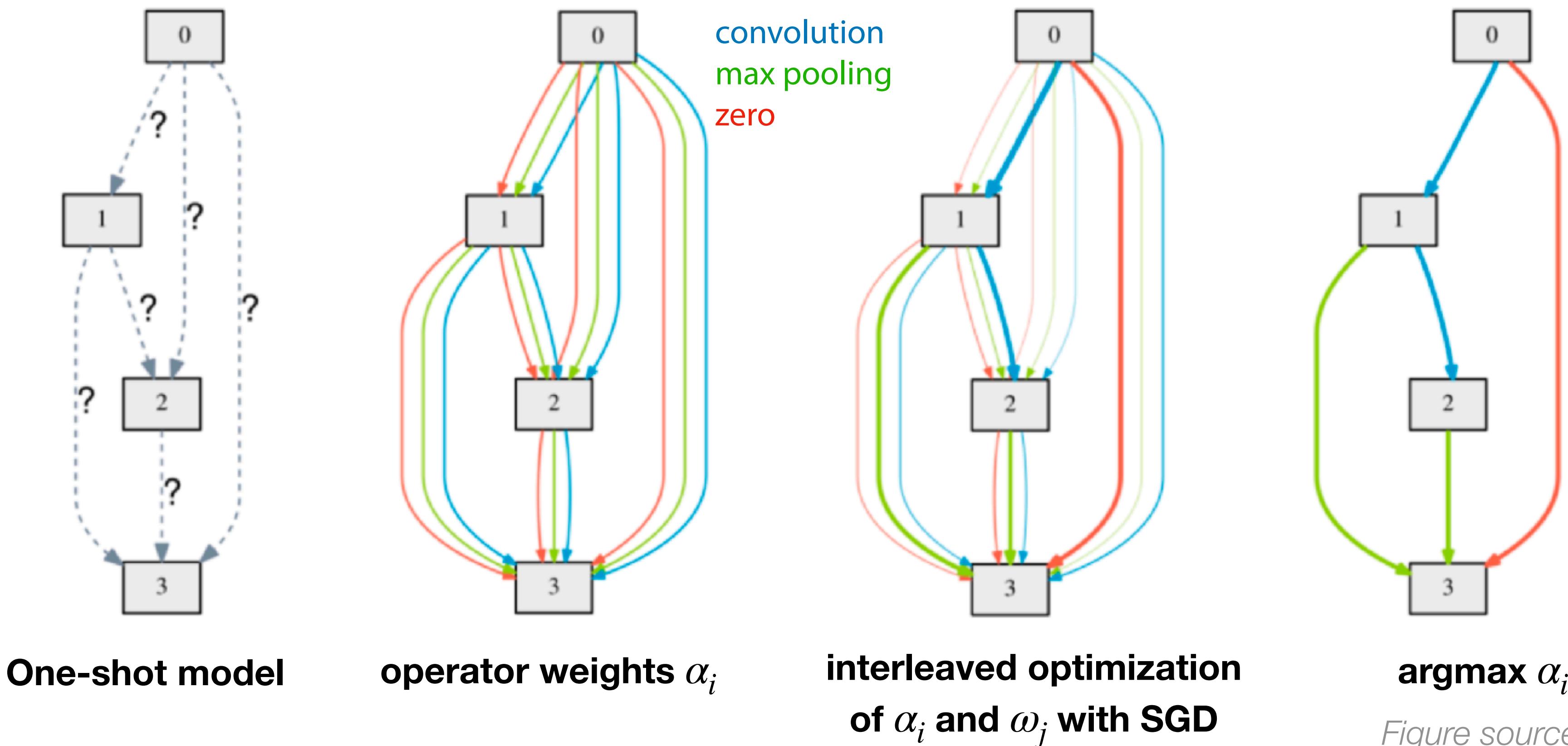
*Figure source: Fusi et al., 2017*

latent representation

$t_{new}$  warm-started  
with  $\lambda_{1..k}$

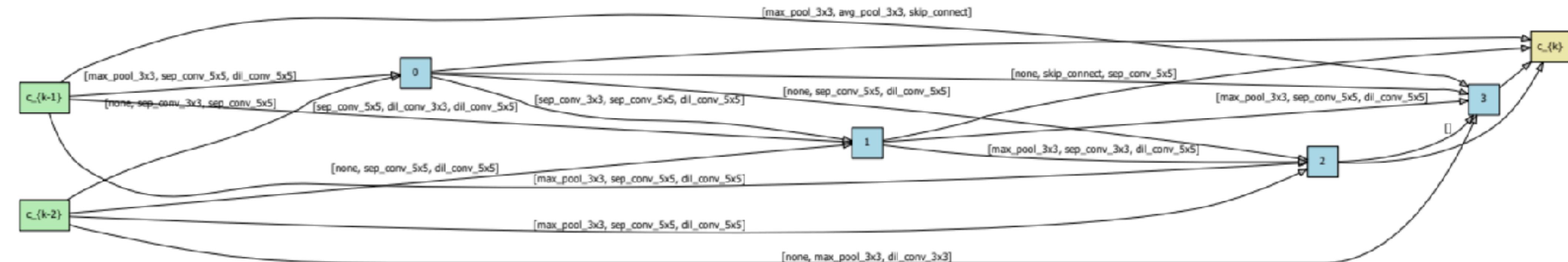
# DARTS: Differentiable NAS

- Fixed (one-shot) structure, learn which operators to use
- Give all operators a weight  $\alpha_i$
- Optimize  $\alpha_i$  and model weights  $\omega_j$  using bilevel optimization
  - approximate  $\omega_j^*(\alpha_i)$  by adapting  $\omega_j$  after every training step

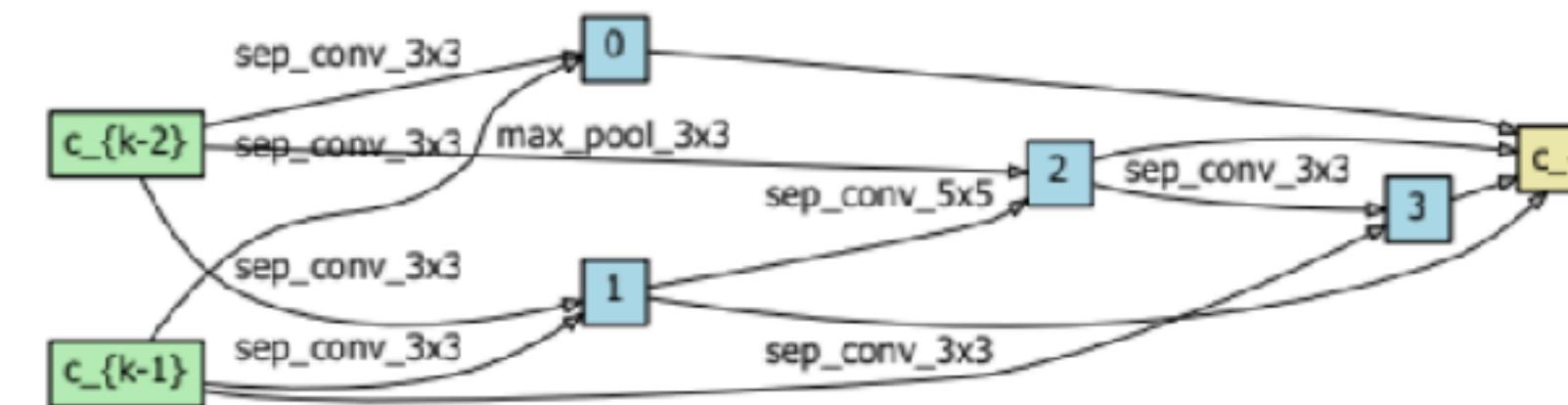
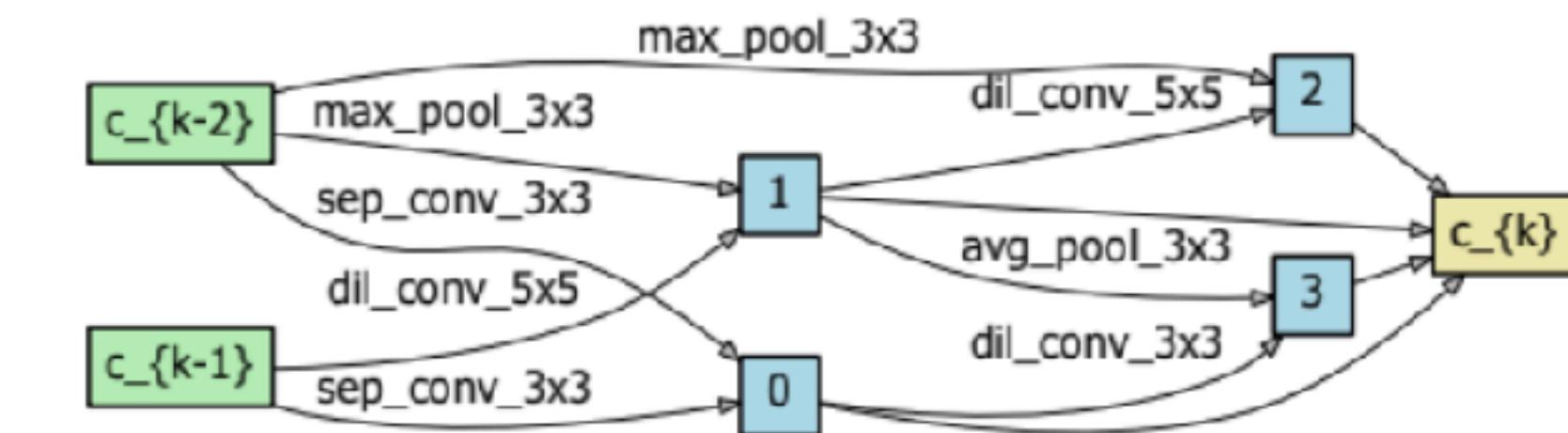


# Warm-started DARTS

- Warm-start DARTS with architectures that worked well on similar problems
- Slightly better performance, but much faster (5x)

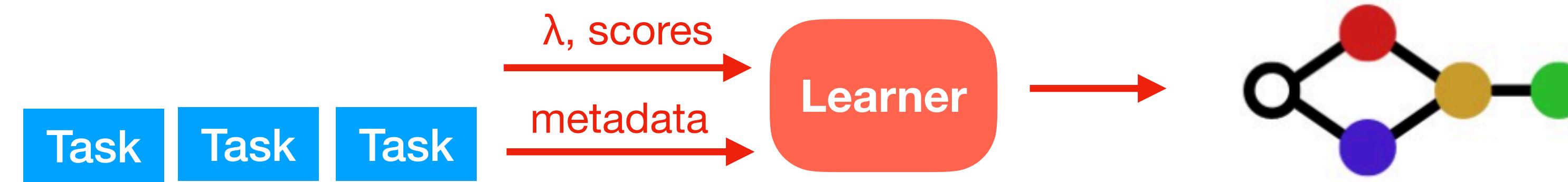


(a) FLOWER\_V3 transfer architecture for normal cell

(b) Normal cell found on *aircraft* task(c) Normal cell found on *birds* task

# Meta-models

(learn how to build models/components)



# Algorithm selection models

- Learn direct mapping between meta-features and  $P_{i,j}$ 
  - Zero-shot meta-models: predict best  $\lambda_i$  given meta-features <sup>1</sup>

$$m_j \rightarrow \text{meta-learner} \rightarrow \lambda_{\text{best}}$$

- Ranking models: return ranking  $\lambda_{1..k}$  <sup>2</sup>

$$m_j \rightarrow \text{meta-learner} \rightarrow \lambda_{1..k}$$

- Predict which algorithms / configurations to consider / tune <sup>3</sup>

$$m_j \rightarrow \text{meta-learner} \rightarrow \Lambda$$

- Predict performance / runtime for given  $\Theta_i$  and task <sup>4</sup>

$$m_j, \lambda_i \rightarrow \text{meta-learner} \rightarrow P_{ij}$$

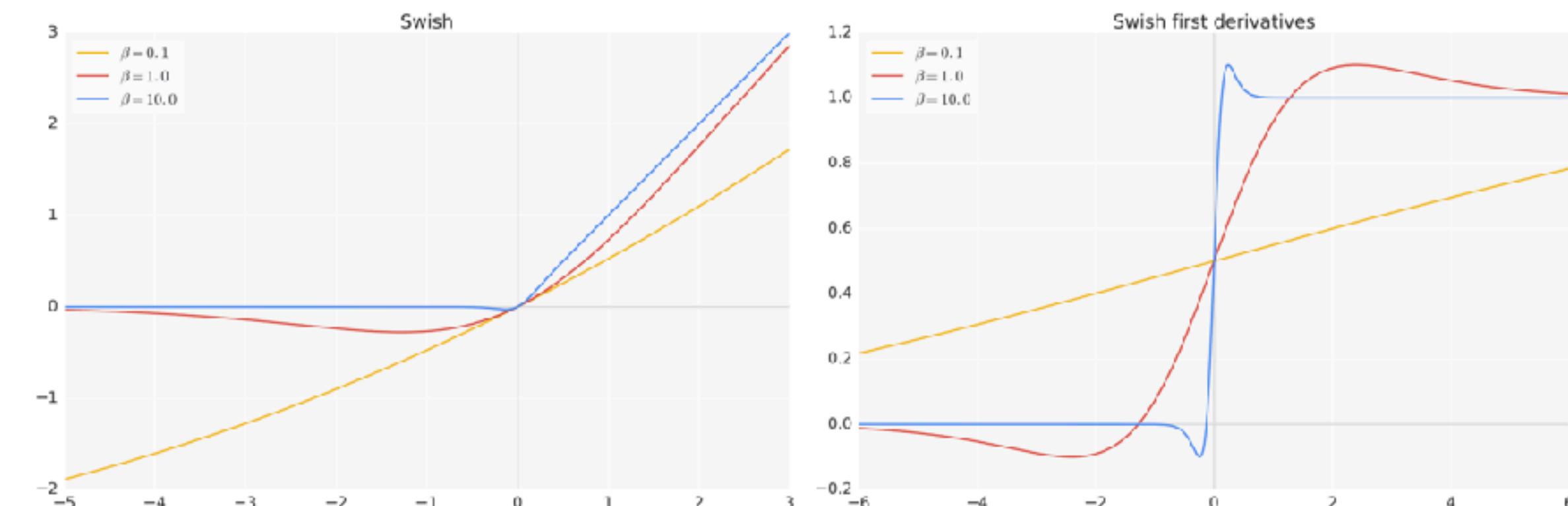
- Can be integrated in larger AutoML systems: warm start, guide search,...

# Learning model components

- Learn nonlinearities: RL-based search of space of likely useful activation functions <sup>1</sup>

- E.g. Swish can outperform ReLU

$$\text{Swish} : \frac{x}{1 + e^{-\beta x}}$$

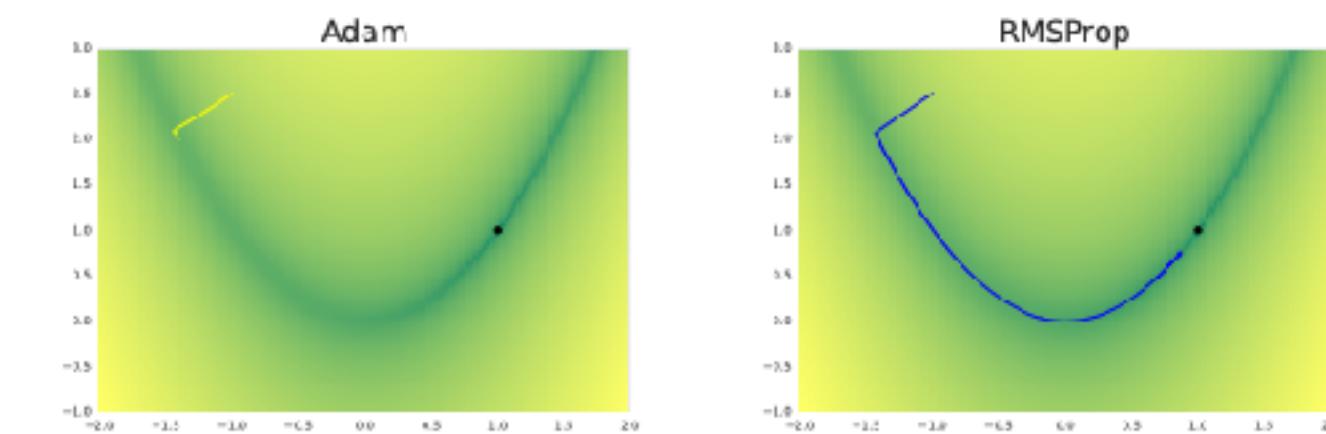


- Learn optimizers: RL-based search of space of likely useful update rules <sup>2</sup>

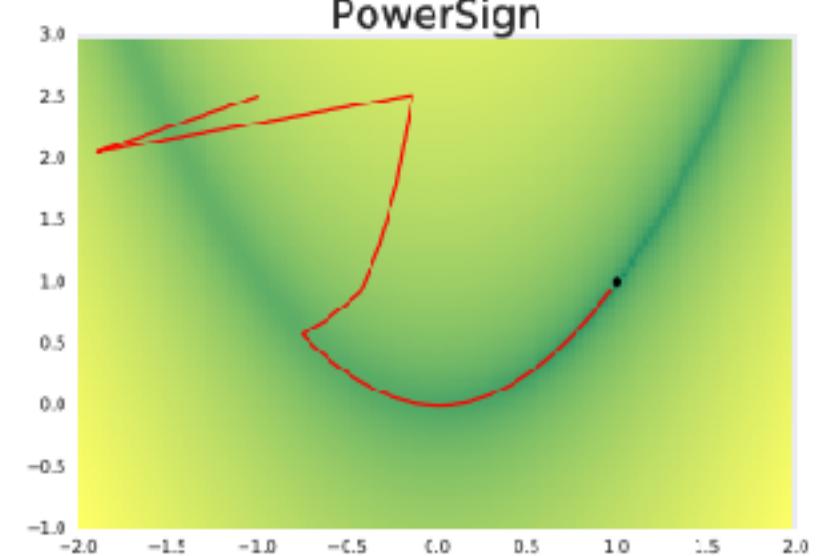
- E.g. PowerSign can outperform Adam, RMSprop

$$\text{PowerSign} : e^{\text{sign}(g)\text{sign}(m)} g$$

*g: gradient, m:moving average*



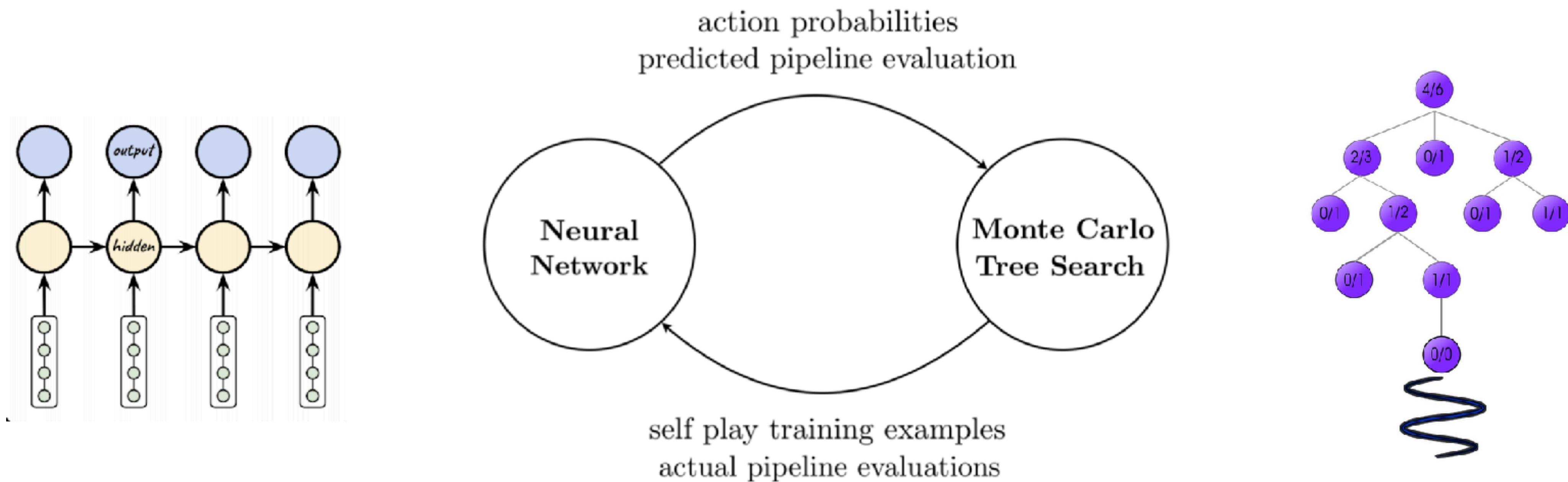
- Learn acquisition functions for Bayesian optimization <sup>3</sup>



# Monte Carlo Tree Search + reinforcement learning

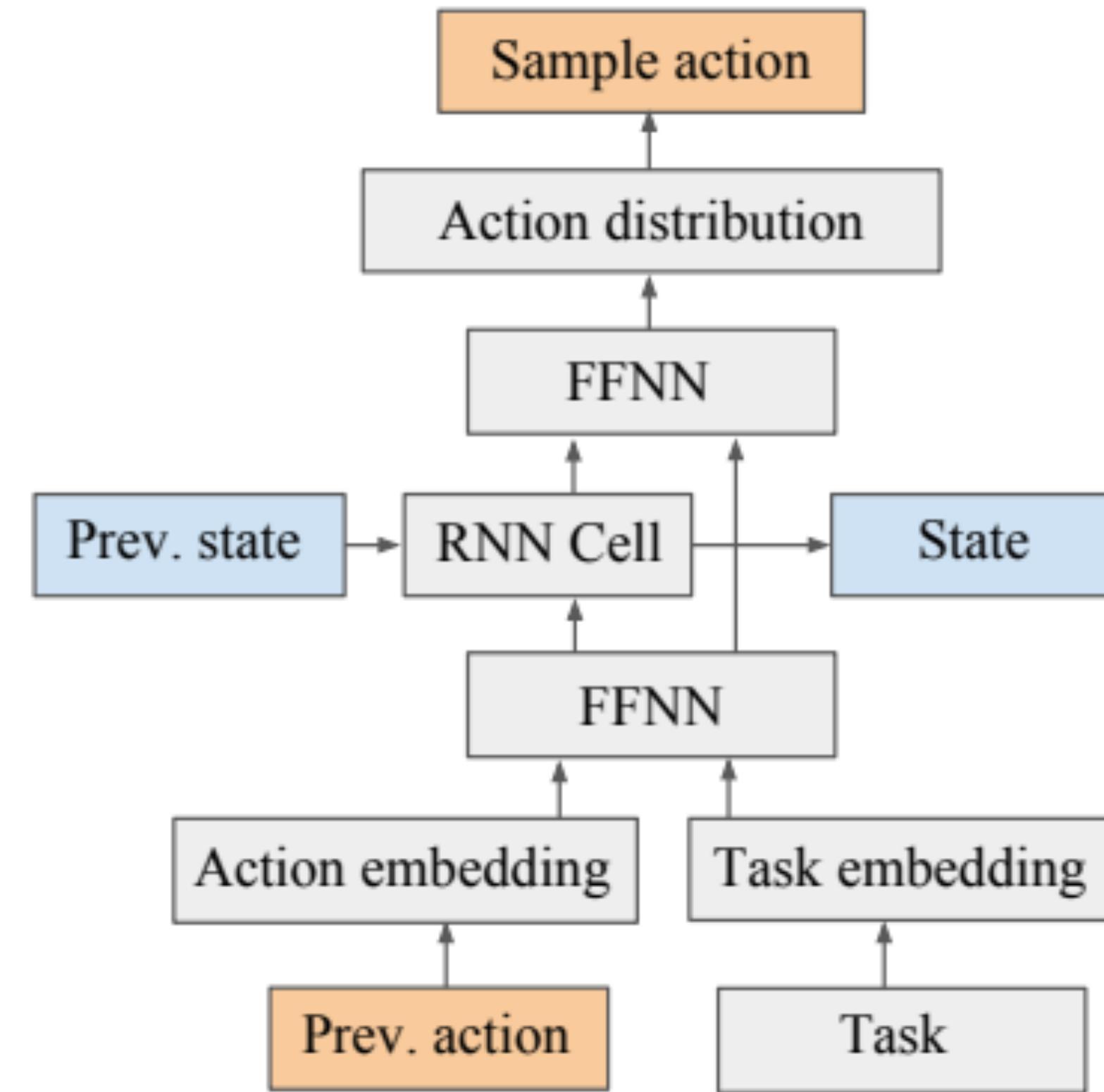
- ***Self-play:***

- Game actions: insert, delete, replace components in a pipeline
- Monte Carlo Tree Search builds pipelines given action probabilities
  - With grammar to avoid invalid pipelines
- Neural network (LSTM) Predicts pipeline performance (can be pre-trained on prior datasets)



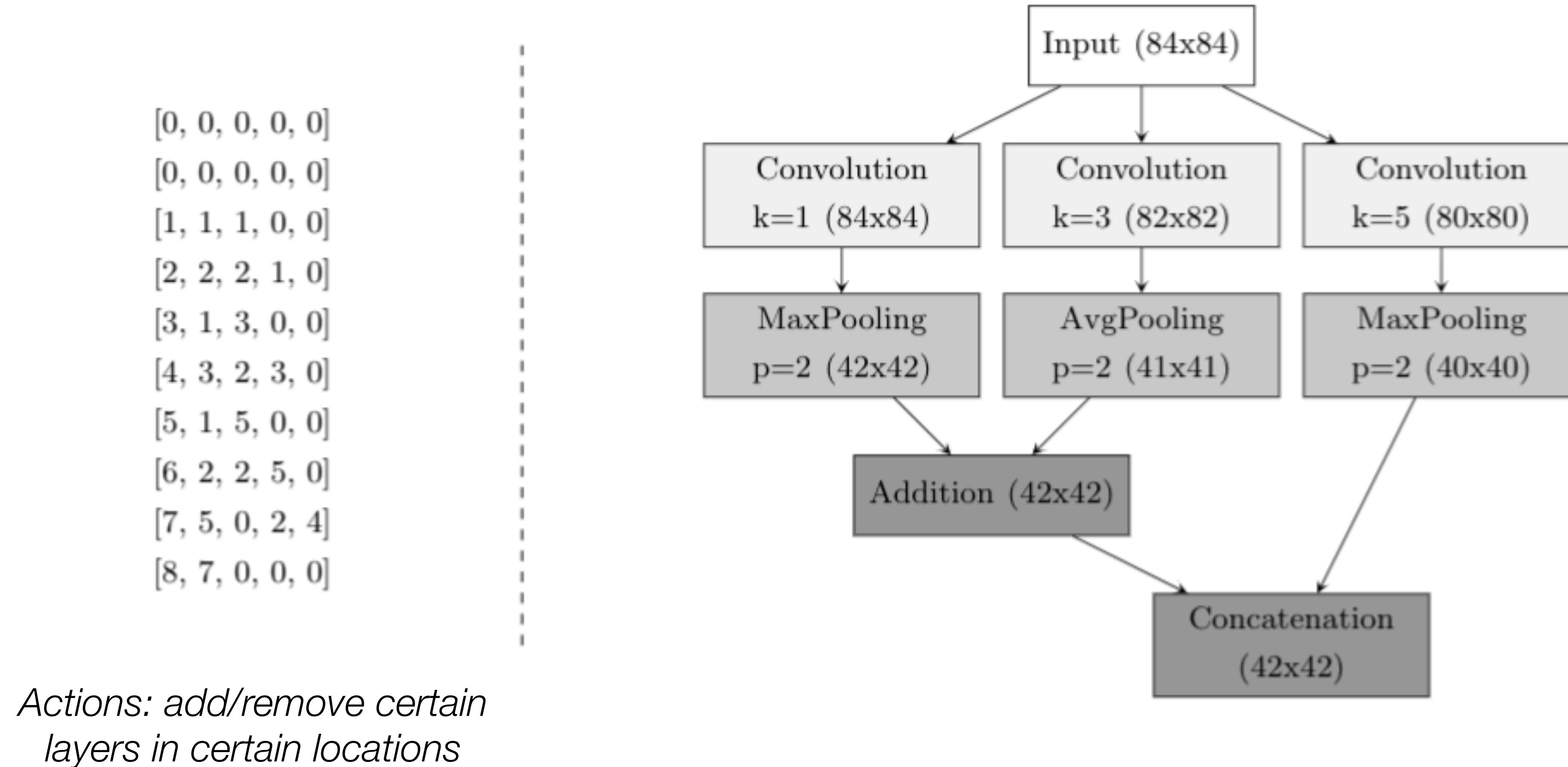
# Neural Architecture Transfer learning

- Warm-start a deep RL controller based on prior tasks
- Much faster than single-task equivalent



# Meta-Reinforcement Learning for NAS

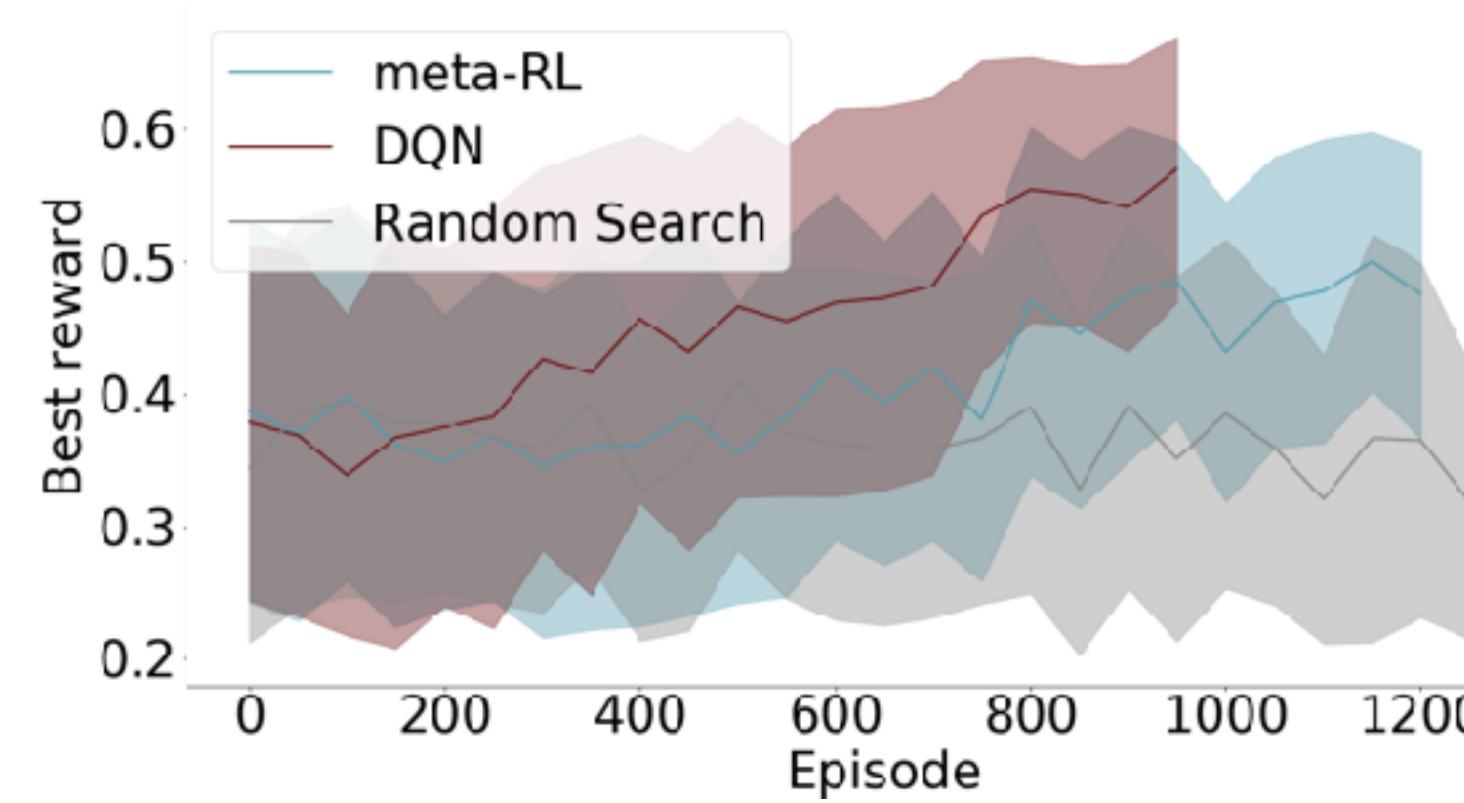
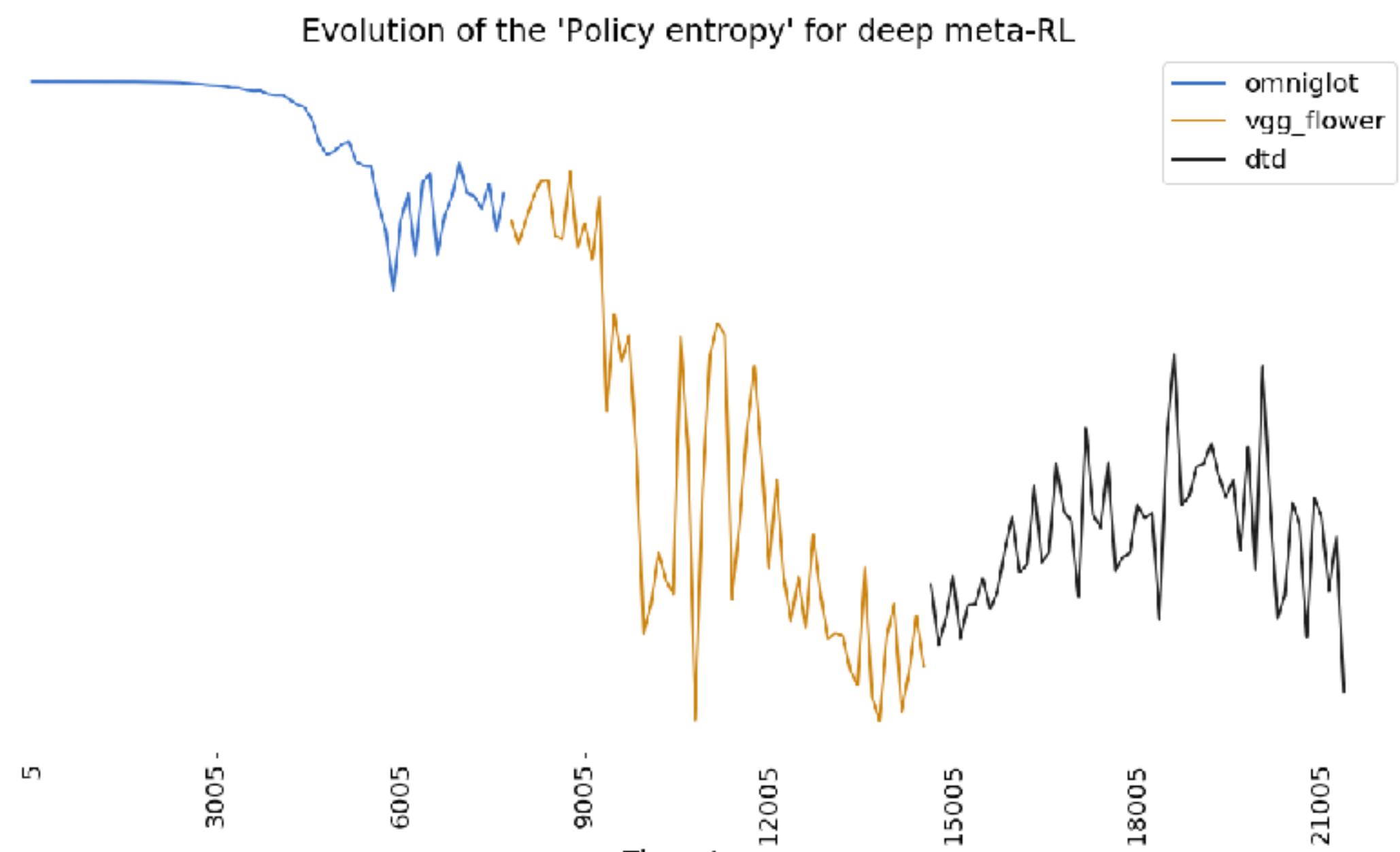
- Train an agent how to build a neural net, across tasks
- Should transfer but also adapt to new tasks



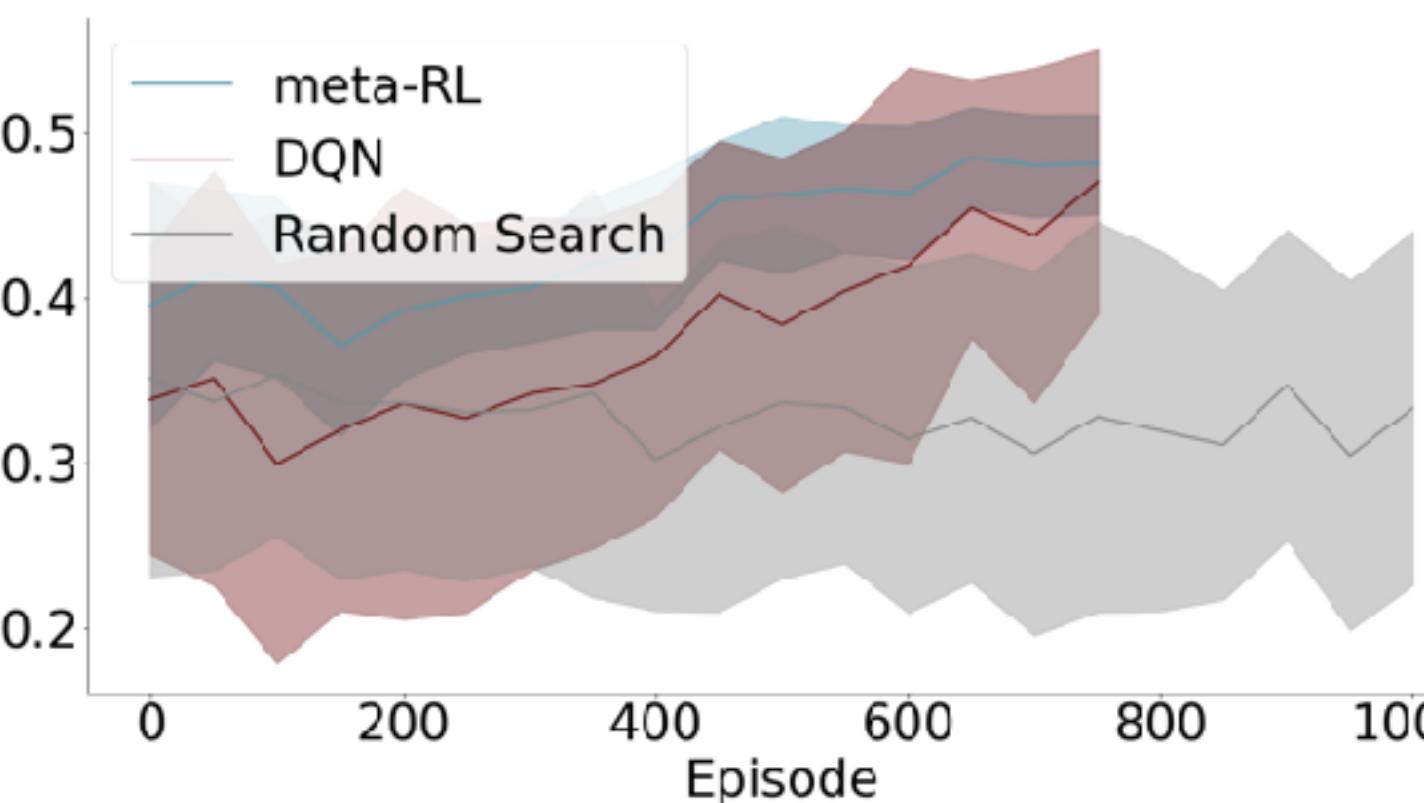
# Meta-Reinforcement Learning for NAS

Results on increasingly difficult tasks:

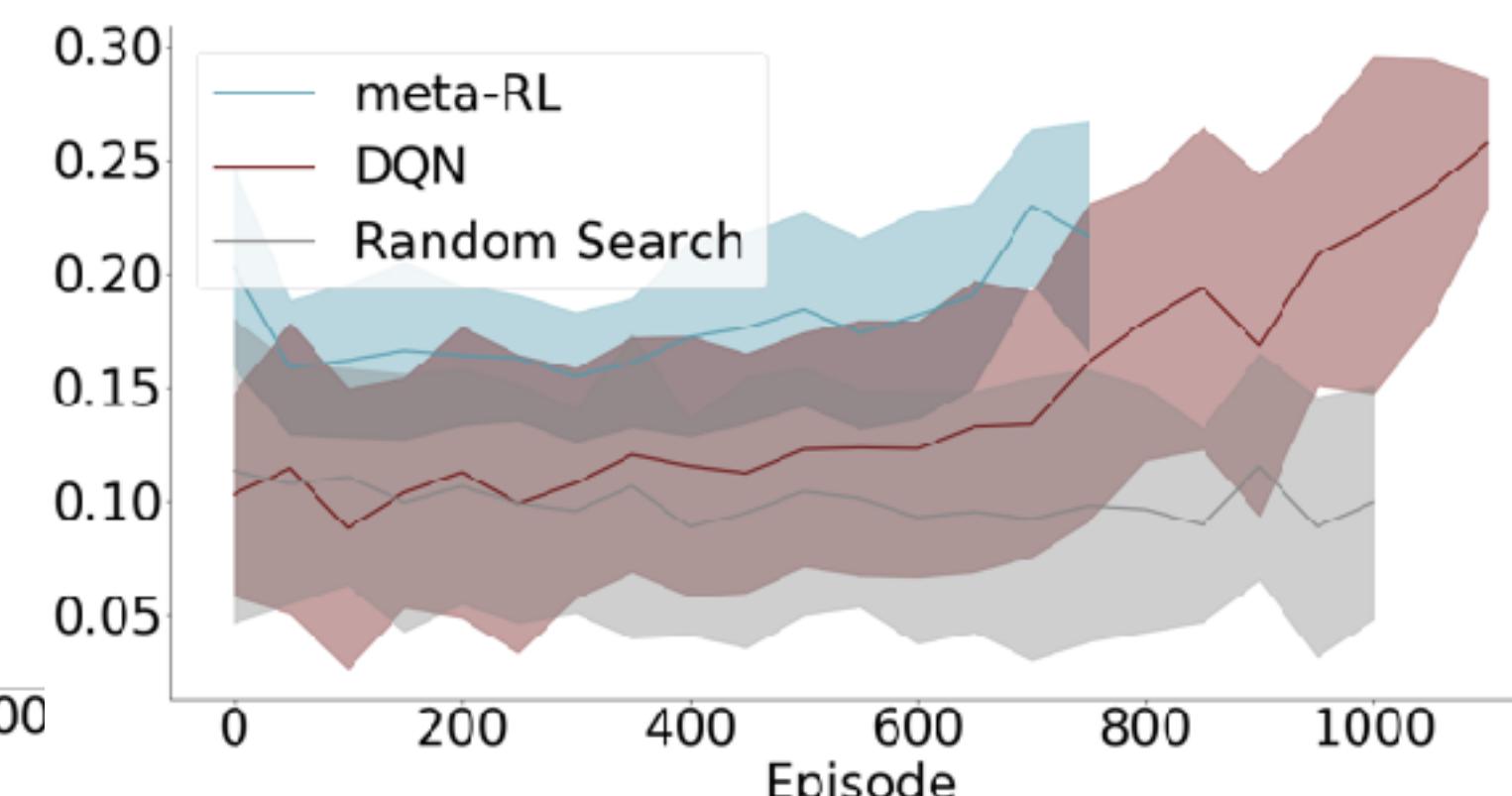
- Initially slower than DQN, but faster after a few tasks
- Policy entropy shows learning/re-learning



omniglot



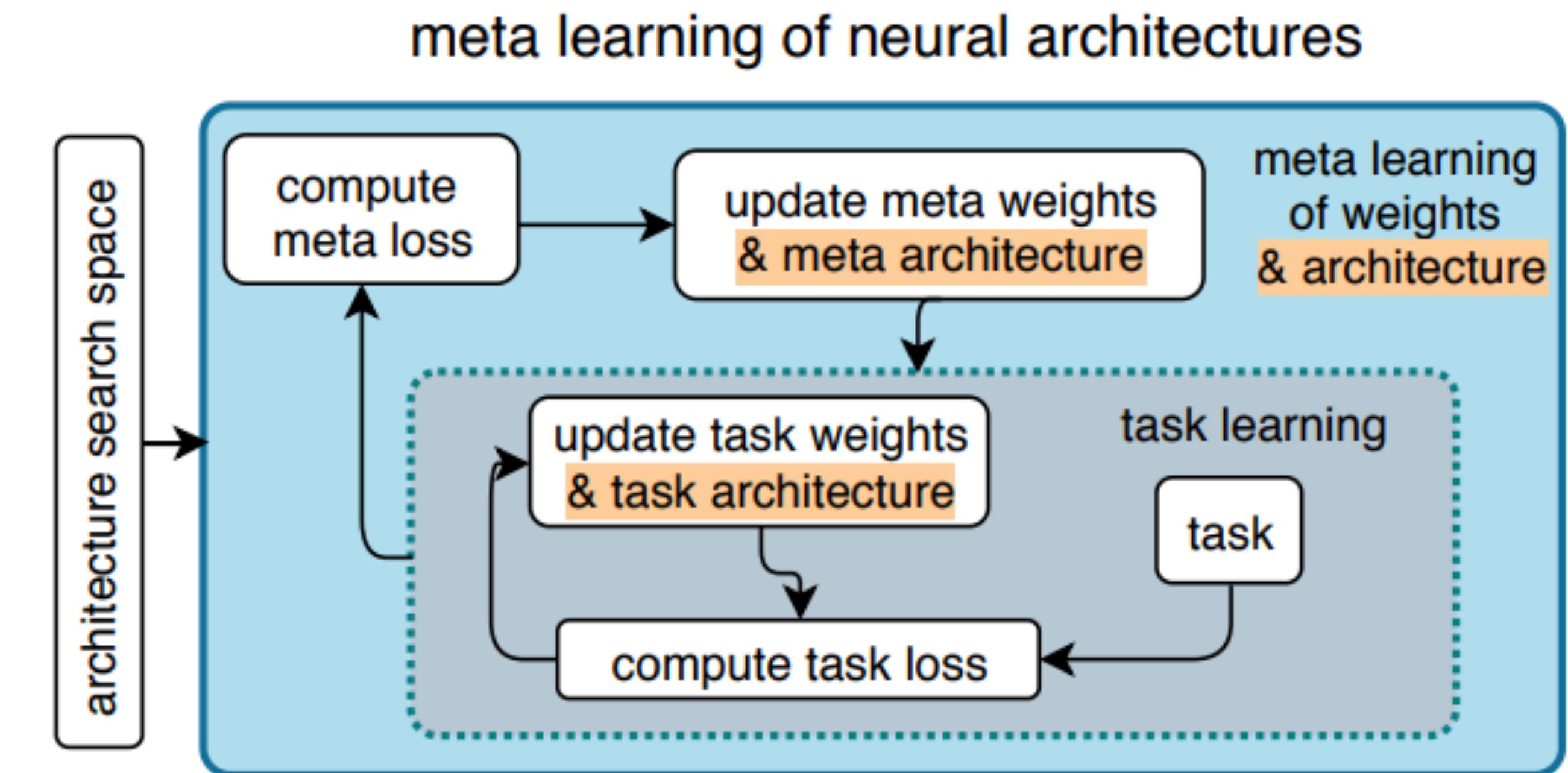
vgg\_flower



dtd

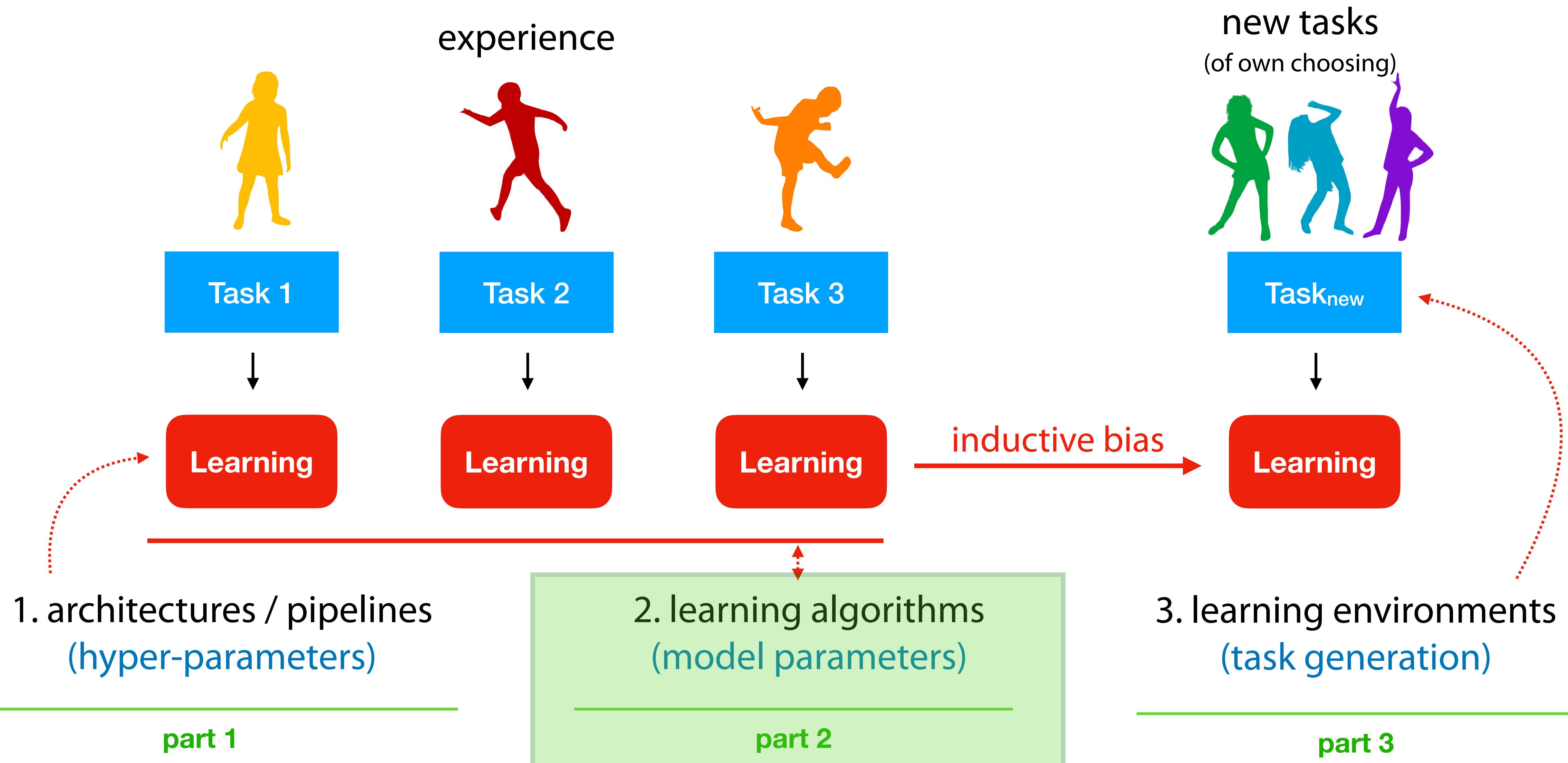
# MetaNAS: MAML + Neural Architecture Search

- Combines gradient based meta-learning (REPTILE) with NAS
- During meta-train, it optimizes the meta-architecture (DARTS weights) along with the meta-parameters (initial weights)  $\Theta$
- During meta-test, the architecture can be adapted to the novel task through gradient descent



# What can we learn to learn?

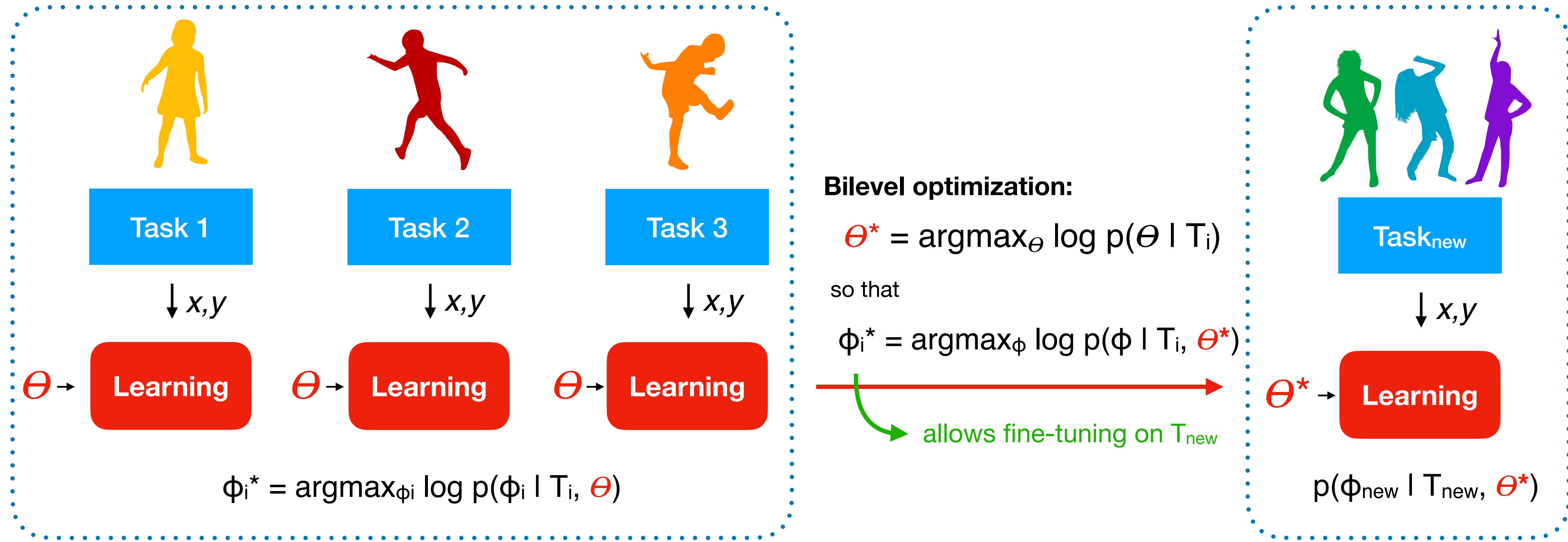
*3 pillars*



# Strategy 1: bilevel optimization

parameterize some aspect of the learner that we want to learn as meta-parameters  $\Theta$

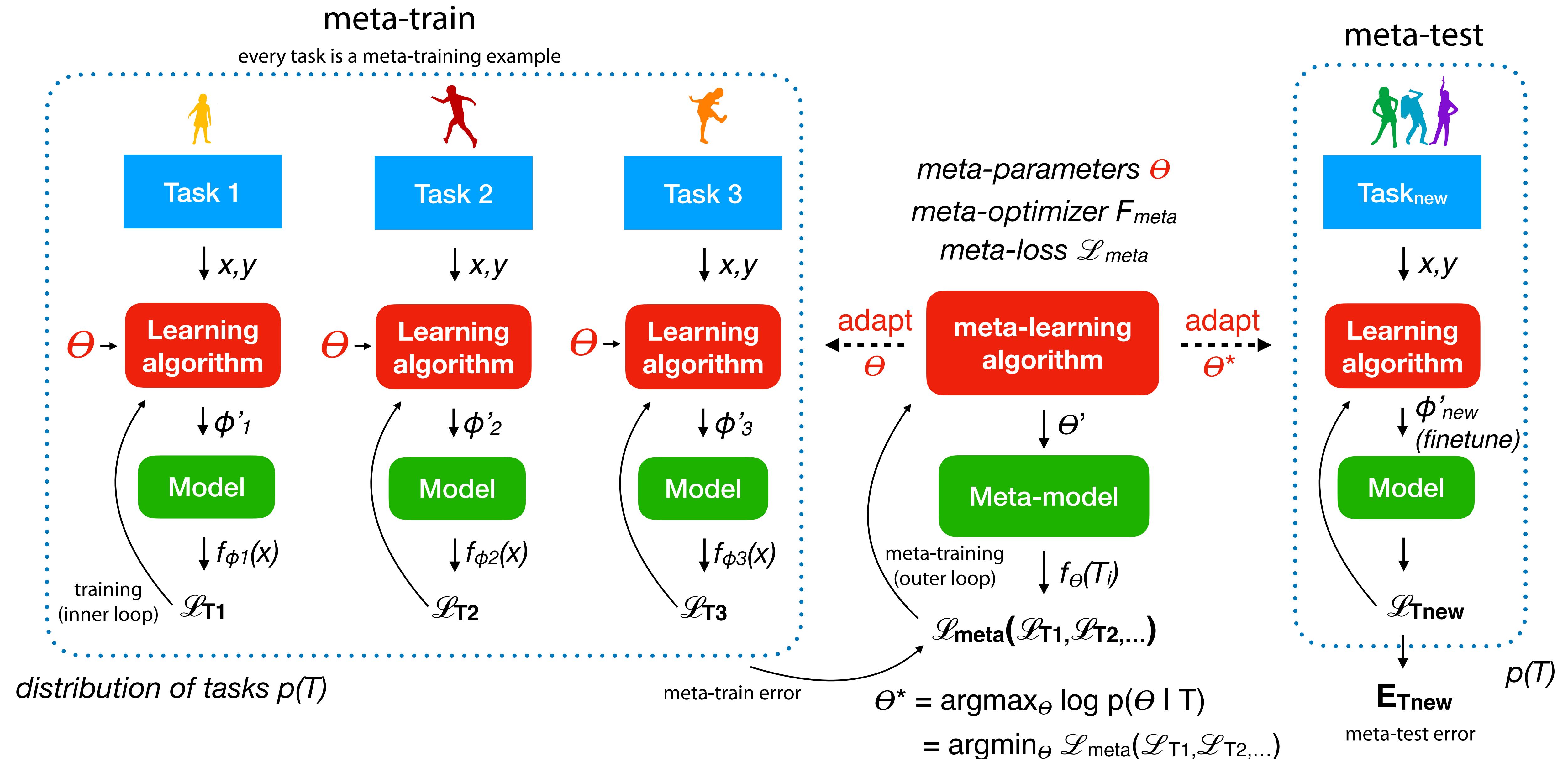
meta-learn  $\Theta$  across tasks



$\Theta$  (prior), could encode an initialization  $\phi$ , the hyperparameters  $\lambda$ , the optimizer,...

Learned  $\Theta^*$  should learn  $T_{new}$  from small amount of data, yet generalize to a large number of tasks

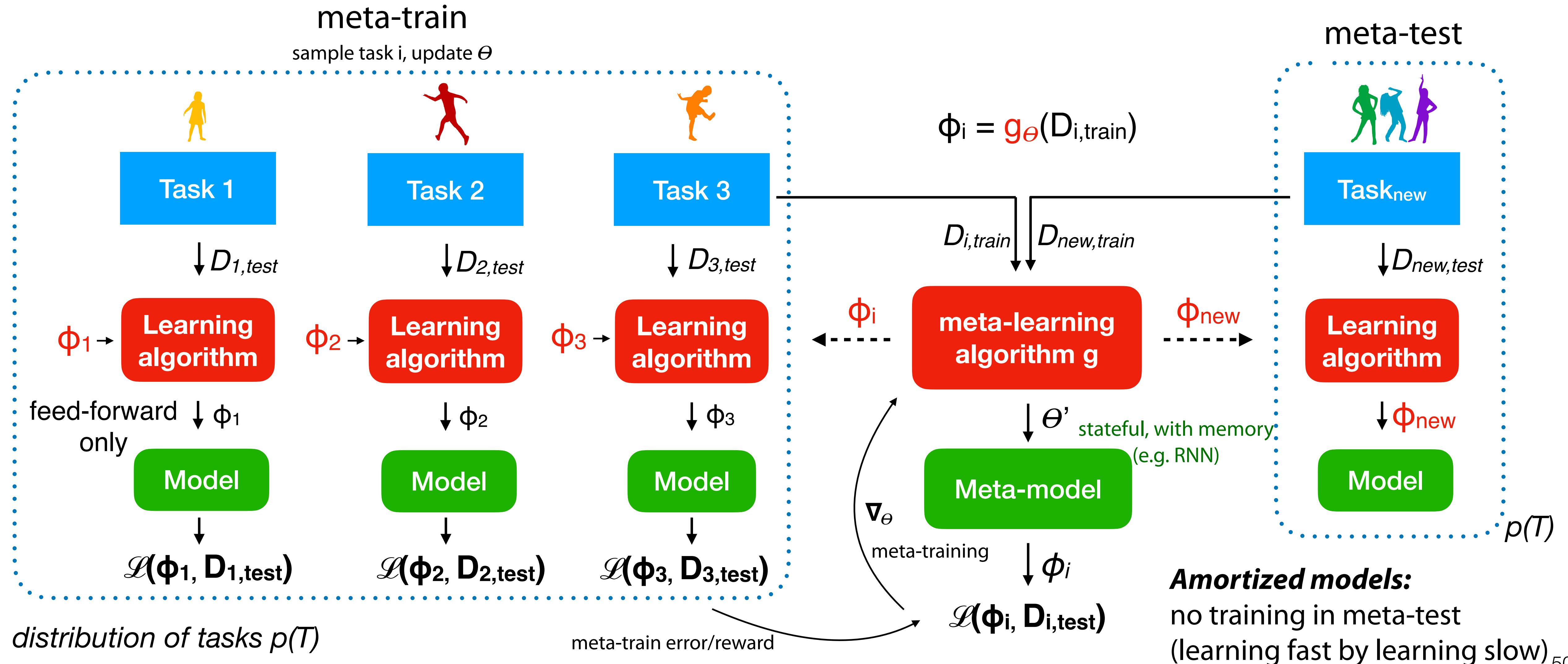
# Meta-learning with bilevel optimization



# Strategy 2: black-box models

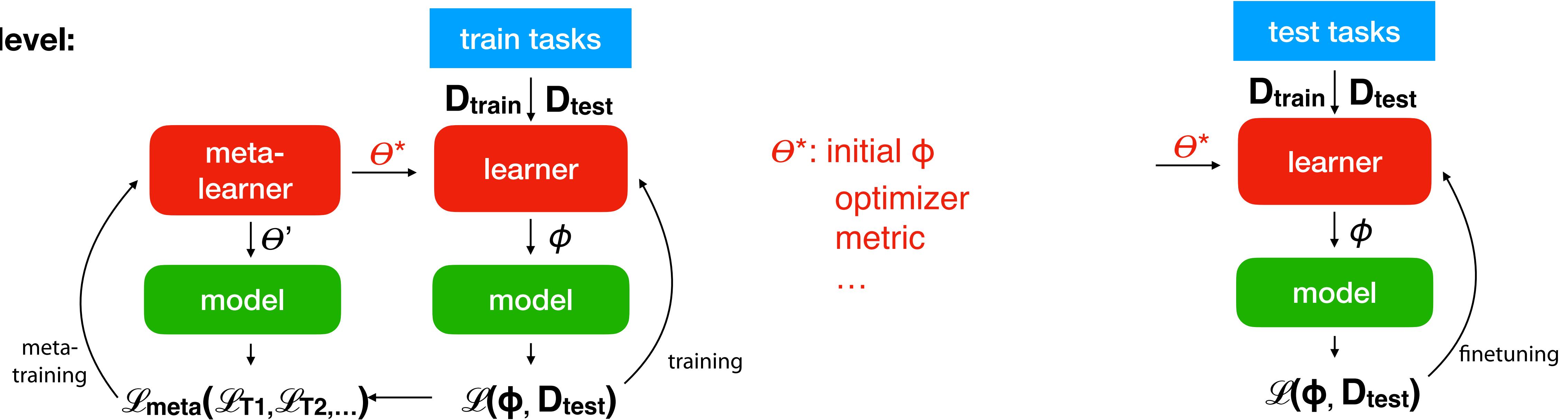
black box meta-model  $g_\theta$ : predicts  $\phi$  given  $D_{train}$  ( $\Theta$  is hidden)

*hypernetwork where input embedding learned across tasks*

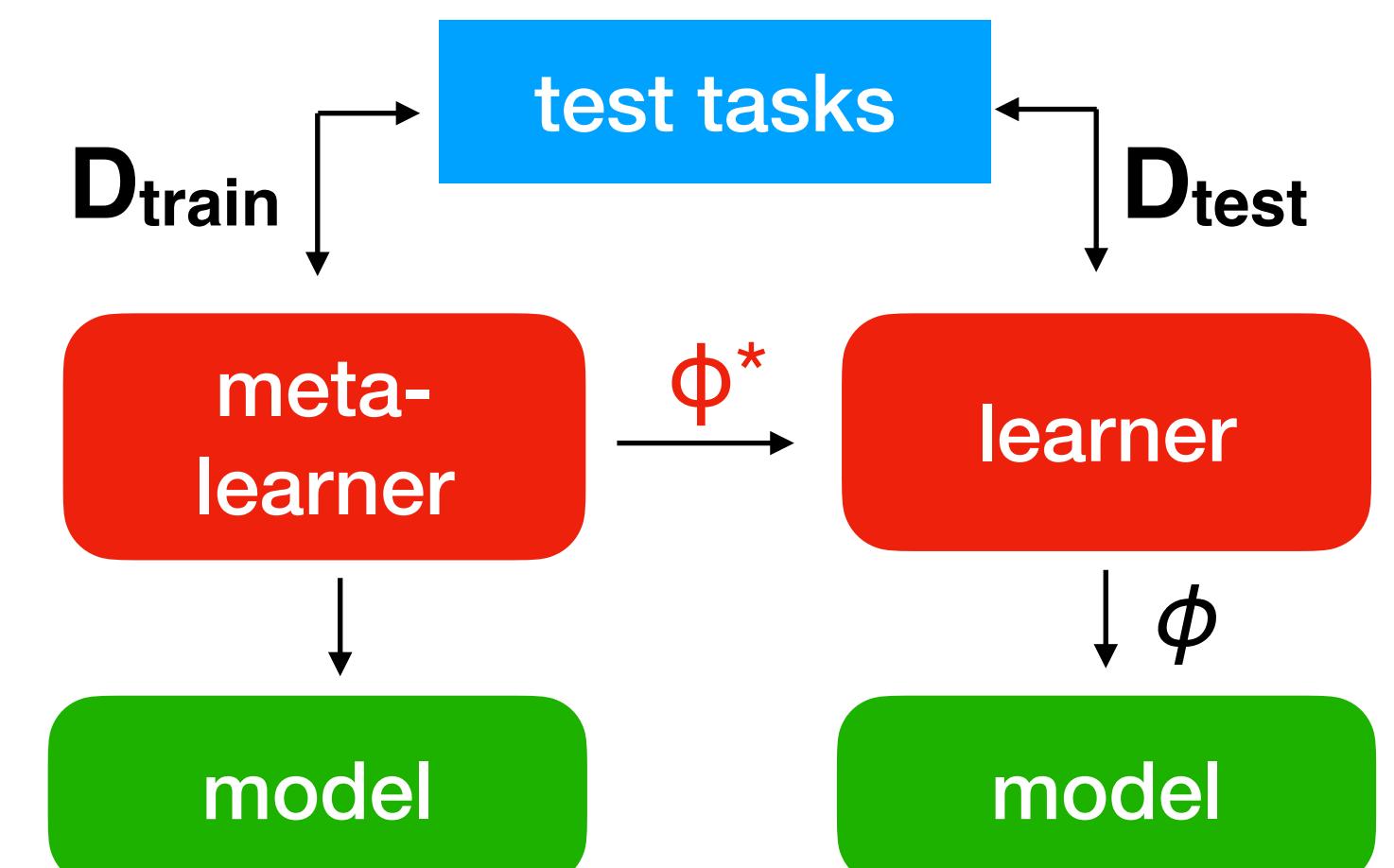
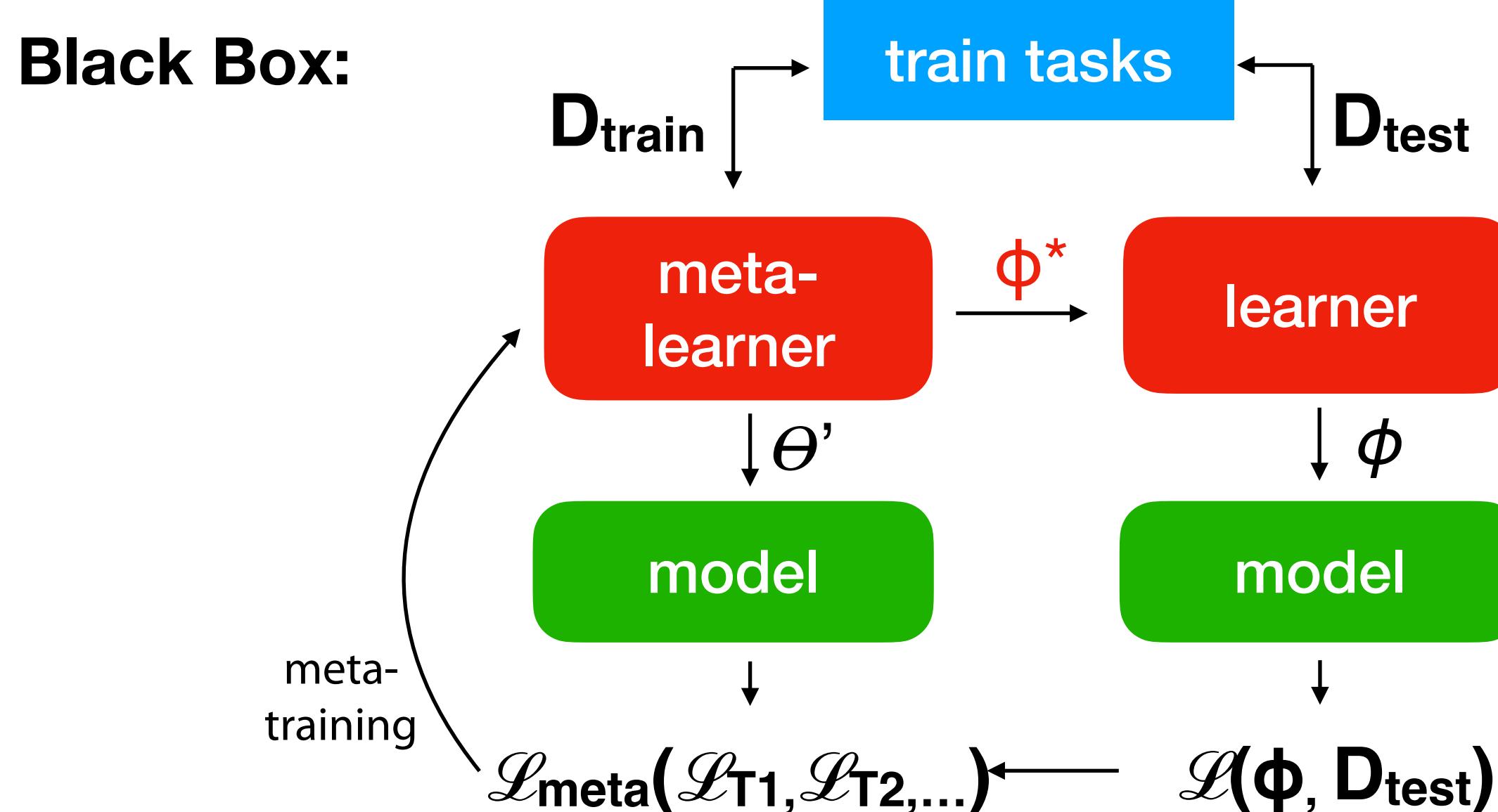


2-way, 3-shot

# Example: few-shot classification

**Bilevel:**

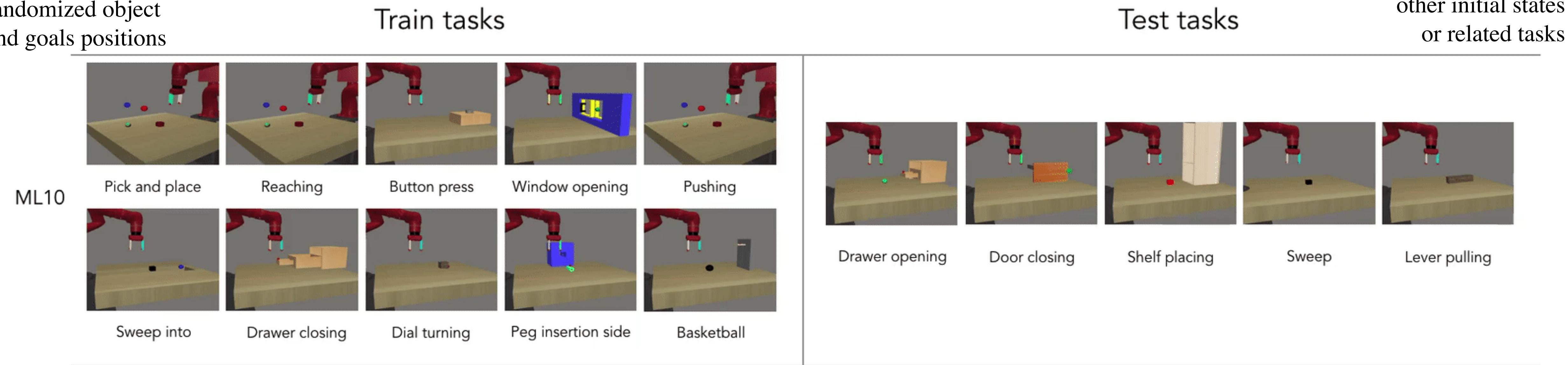
# Example: few-shot classification



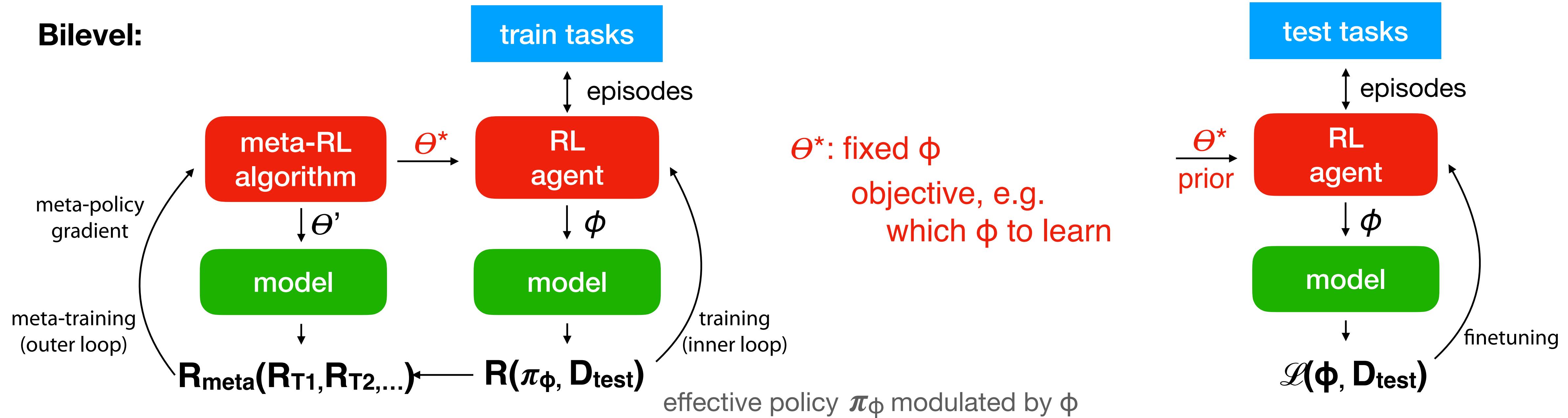
# Example: meta-reinforcement learning

initial state:

randomized object  
and goals positions



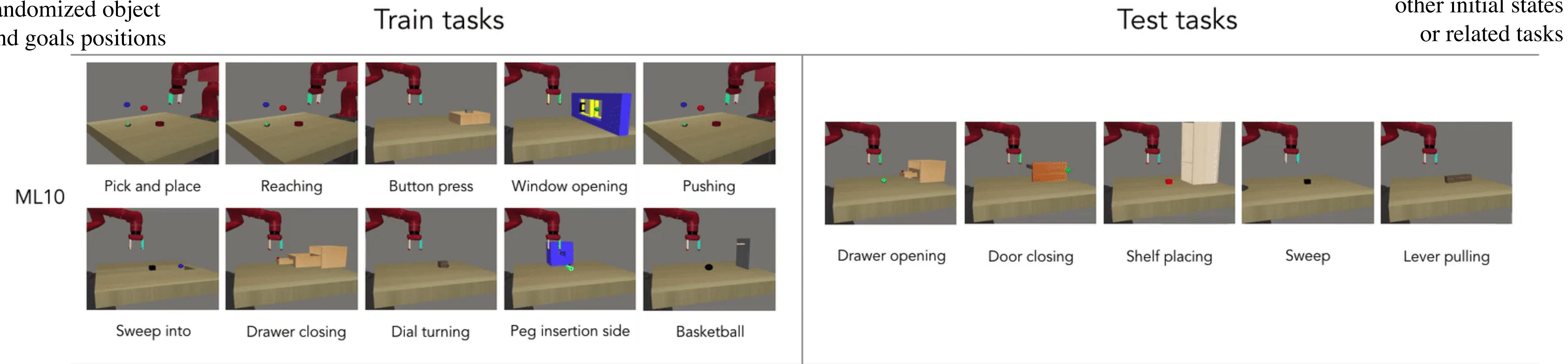
**Bilevel:**



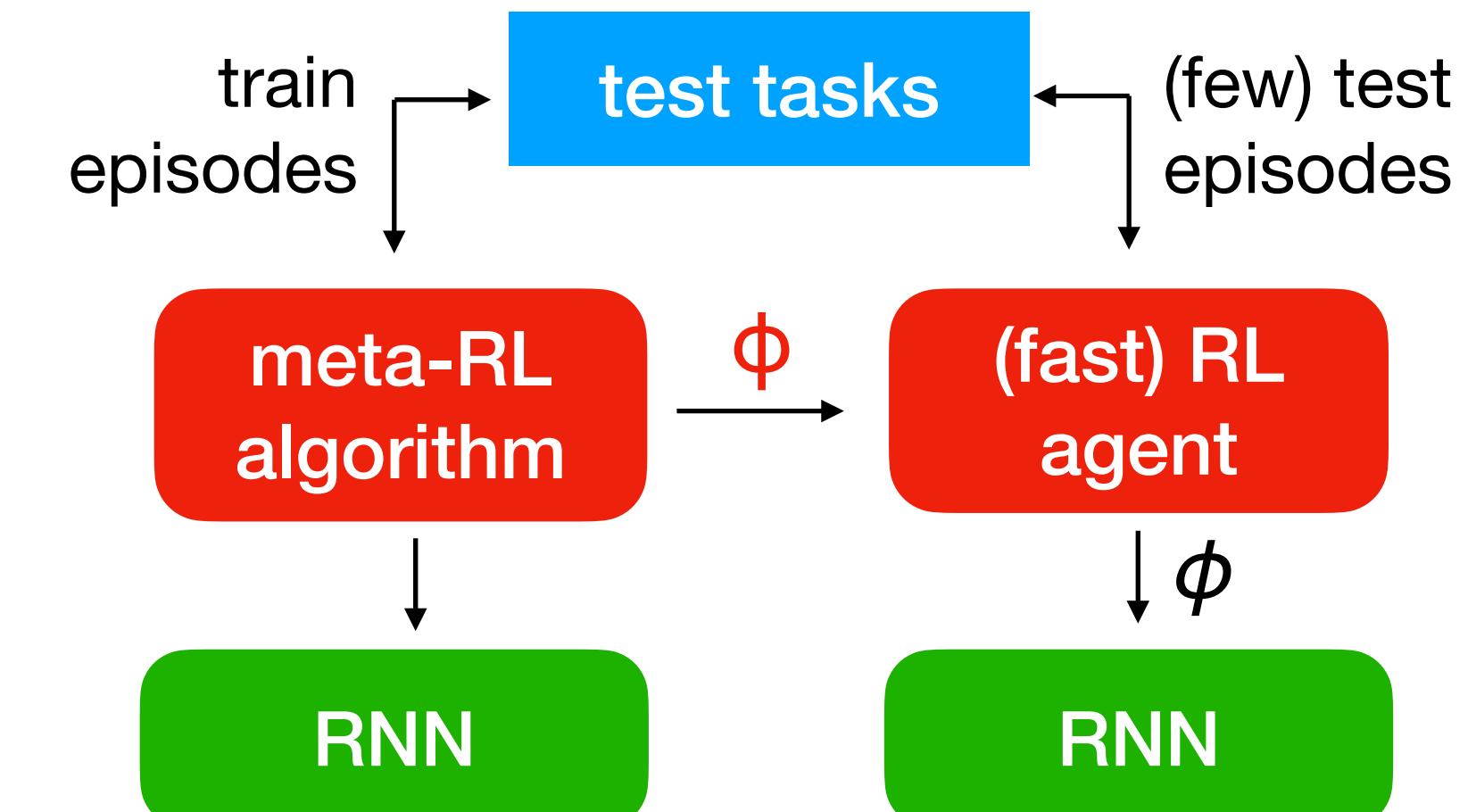
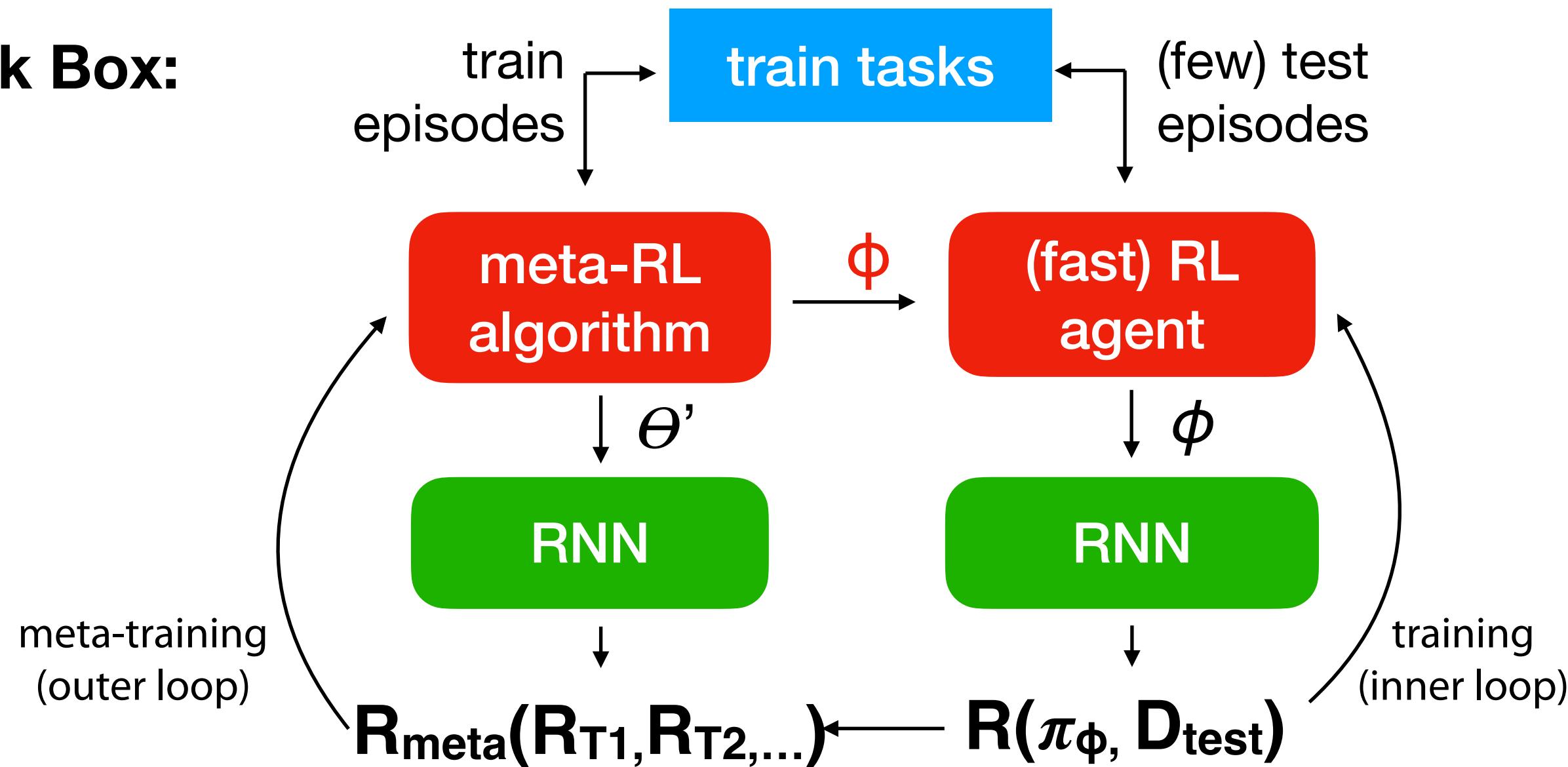
# Example: meta-reinforcement learning

initial state:

randomized object  
and goals positions

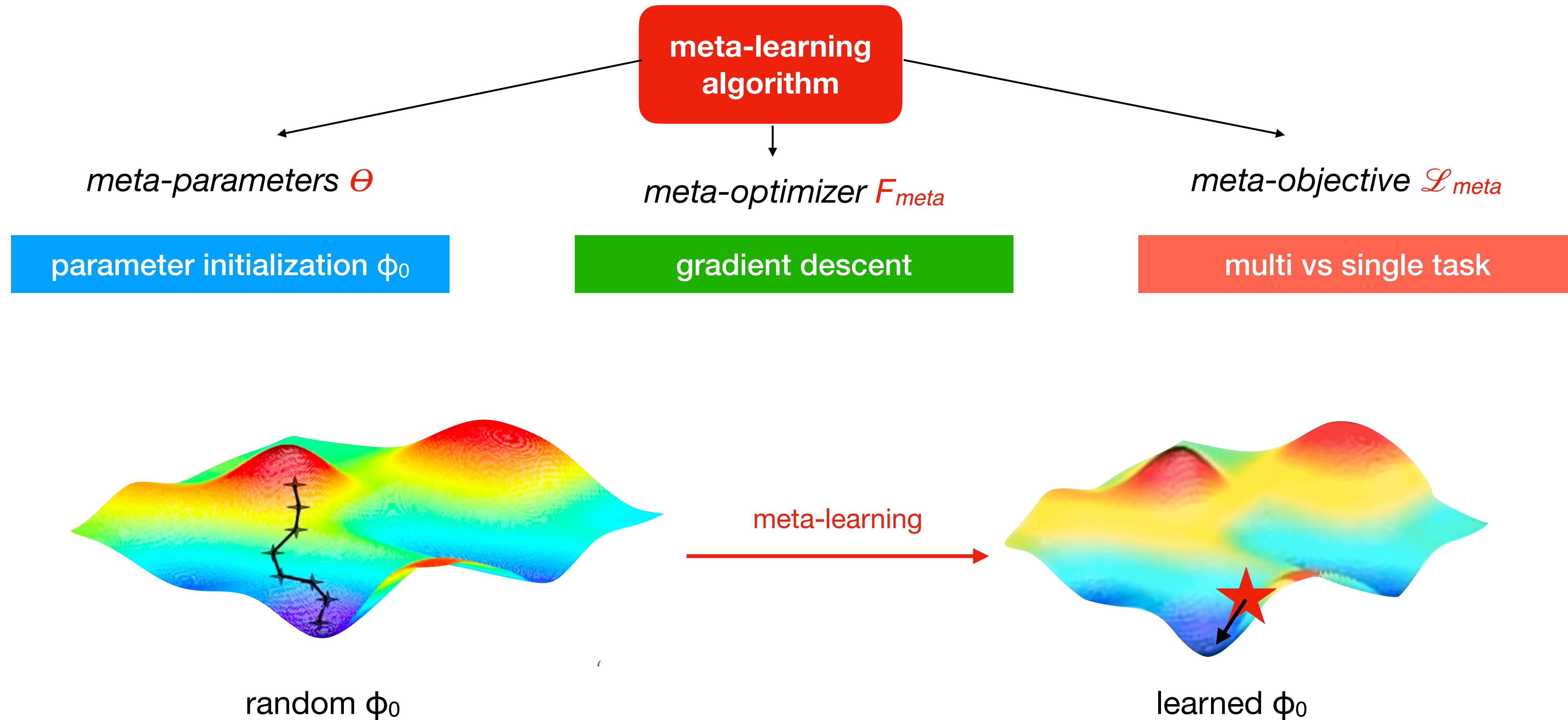


**Black Box:**



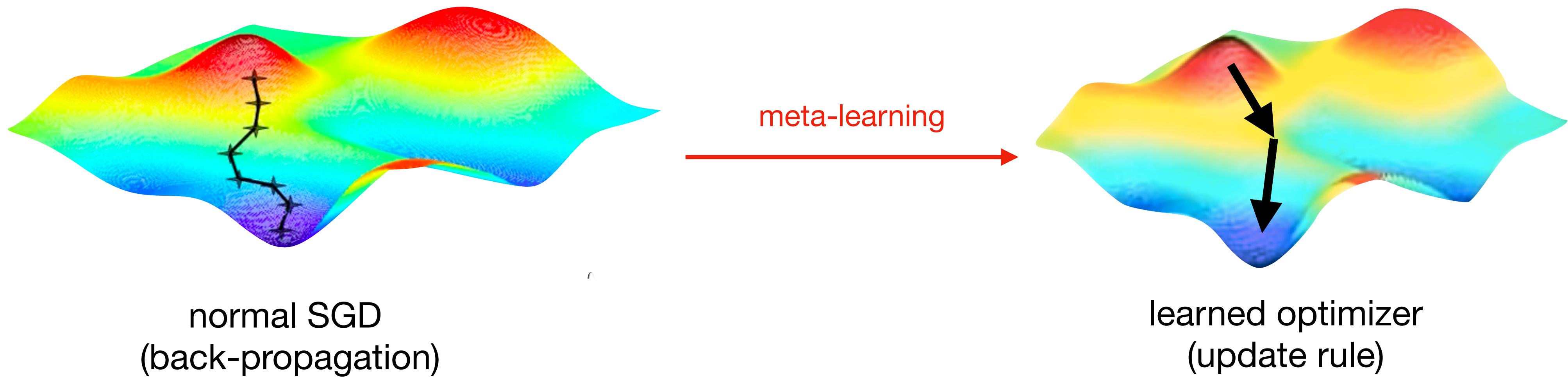
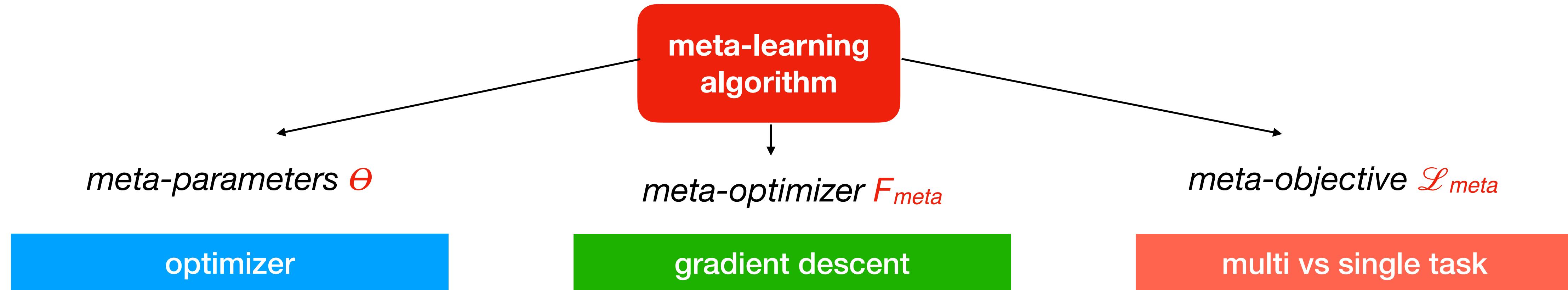
# Taxonomy of meta-learning methods

like base-learners, meta-learners consist of a representation, an objective, and an optimizer



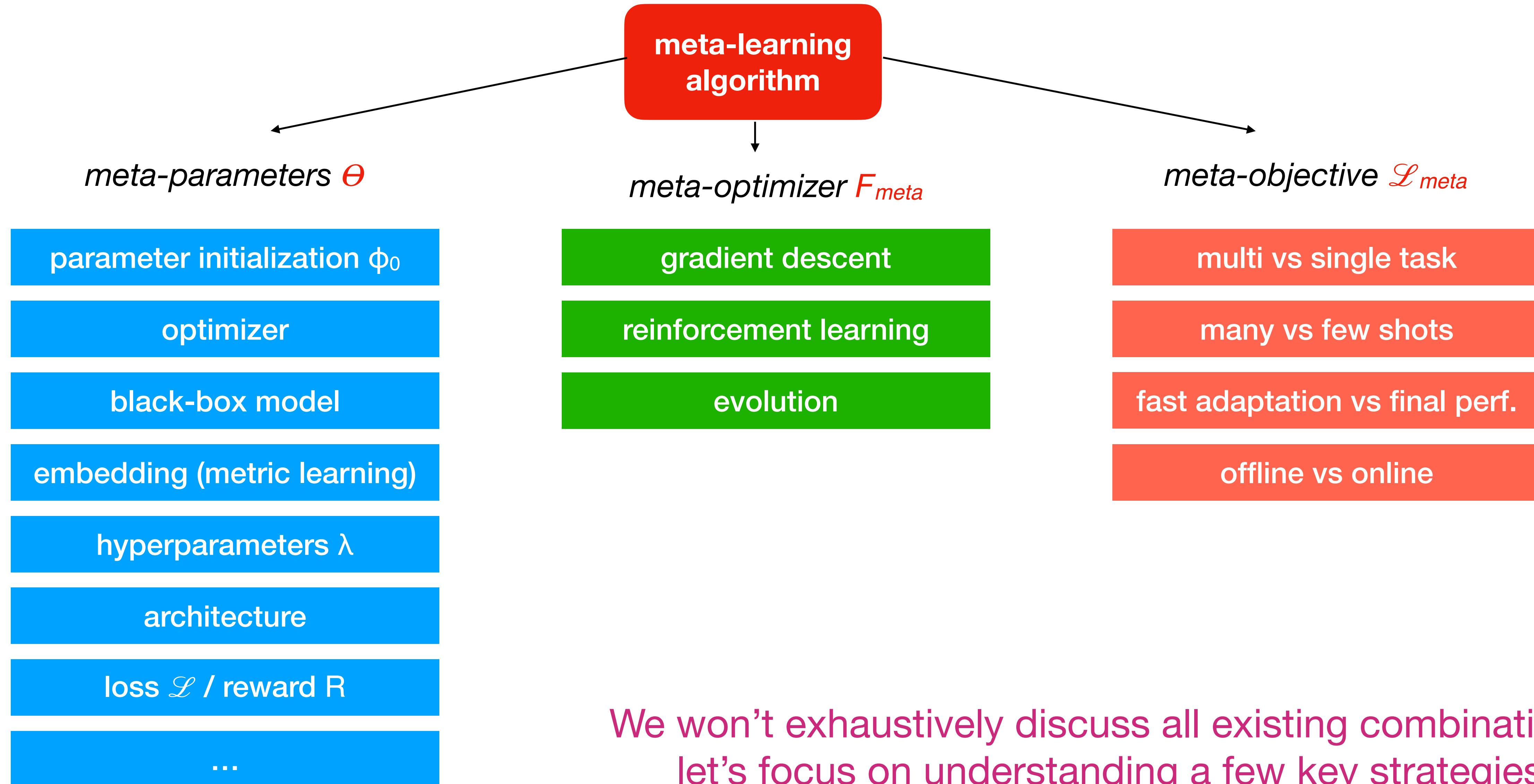
# Taxonomy of meta-learning methods

like base-learners, meta-learners consist of a representation, an objective, and an optimizer

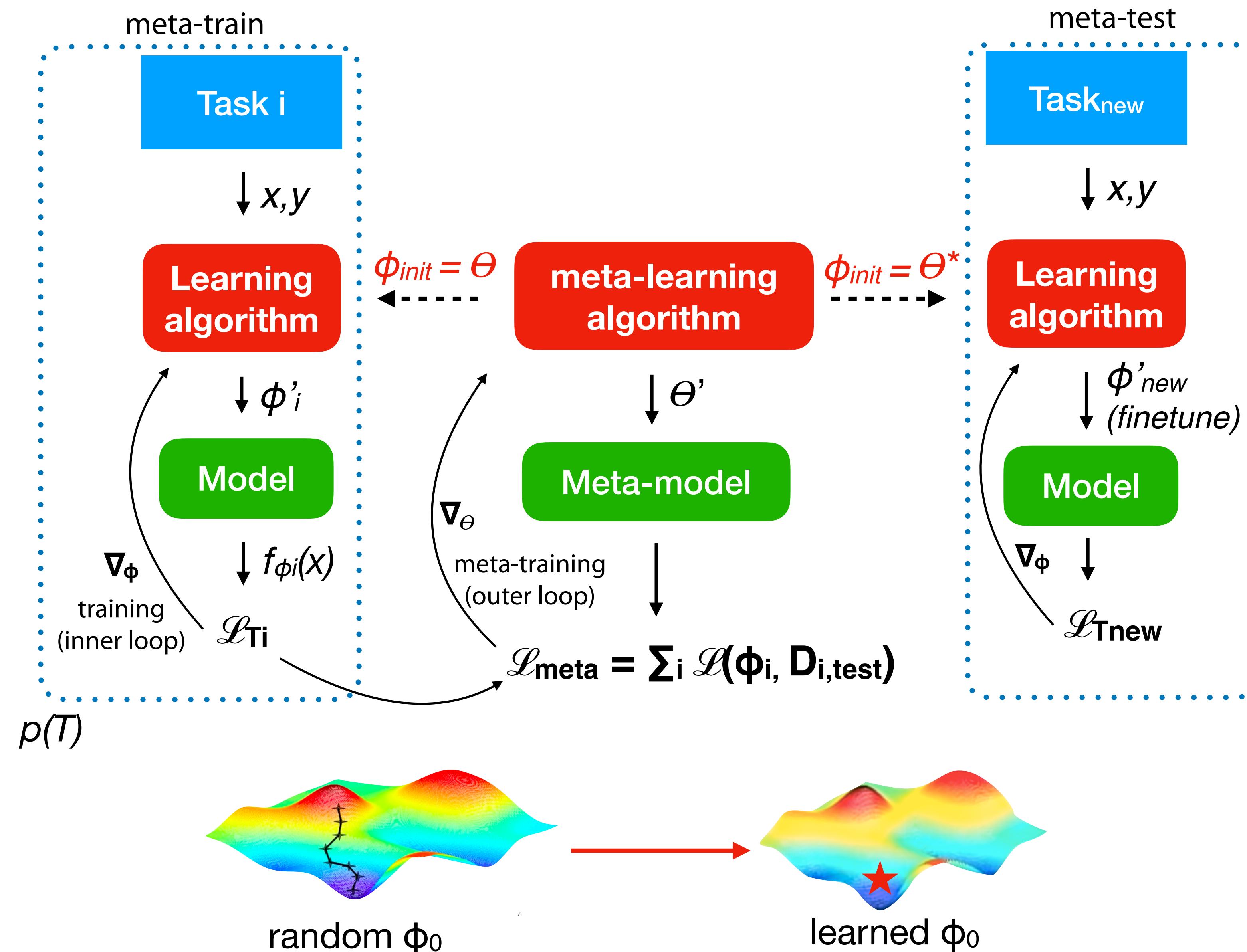


# Taxonomy of meta-learning methods

like base-learners, meta-learners consist of a representation, an objective, and an optimizer



# Gradient-based methods: learning $\phi_{\text{init}}$



- $\Theta$  (prior): model initialization  $\phi_{\text{init}}$
- learn representation suitable for many tasks (e.g. pretrained CNN)
- maximize rapid learning
- Each task  $i$  yields task-adapted  $\phi_i$ 
  - Update algorithm  $u$   
 $\phi_i = u(\Theta, D_{i,\text{train}})$
- Finetune  $\Theta^*$  on  $T_{\text{new}}$  (in few steps)  
 $\phi'^{\text{new}} = u(\Theta^*, D_{\text{new},\text{train}})$

Can be seen as bilevel optimization:

$$\Theta^* = \operatorname{argmin}_{\Theta} \sum_i \mathcal{L}_i(\phi_i, D_{i,\text{test}})$$

$$\phi_i = u(\Theta, D_{i,\text{train}})$$

# Model agnostic meta-learning (MAML)

## Meta-training

- Current initialization  $\Theta$ , model  $f_\Theta$
- On  $i$  tasks, perform  $k$  SGD steps to find  $\phi_i^*$ , then evaluate  $\nabla_\theta \mathcal{L}_i(f_{\phi_i^*})$
- Update task-specific parameters:  $\phi_i = \theta - \alpha \nabla_\theta \mathcal{L}_i(f_{\phi_i^*})$
- Update  $\Theta$  to minimize sum of per-task losses, repeat

$$\theta \leftarrow \theta - \beta \nabla_\theta \sum_i \mathcal{L}_i(f_{\phi_i})$$

$\alpha, \beta$ : learning rates

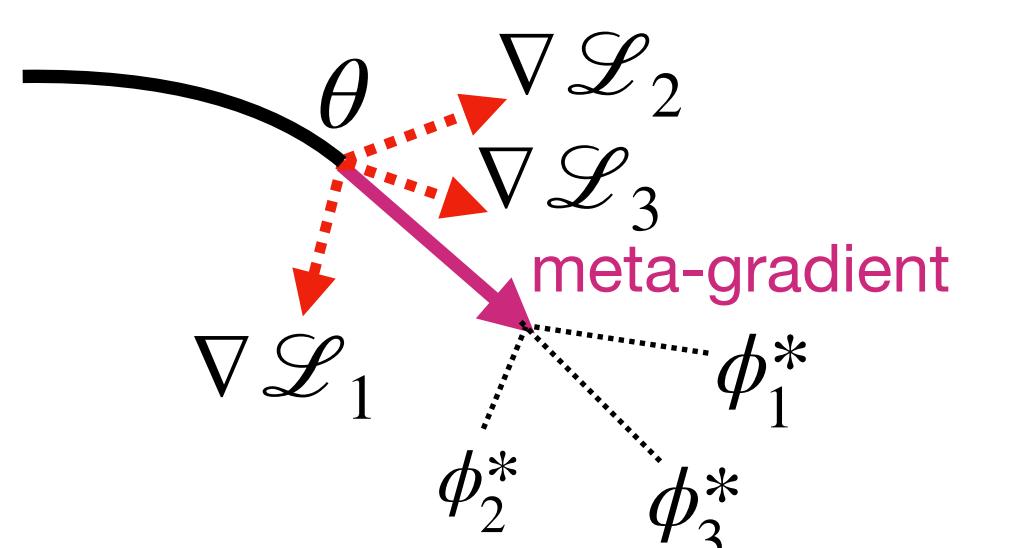
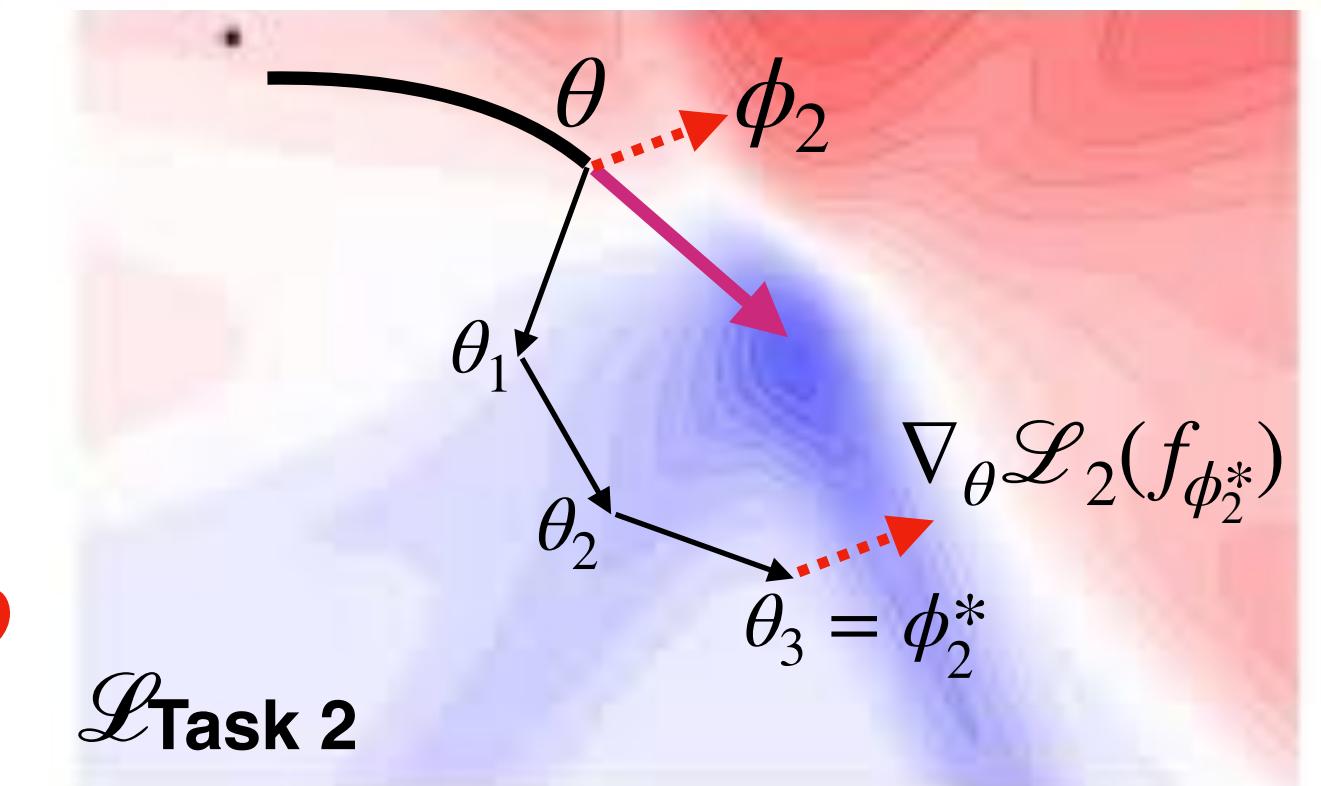
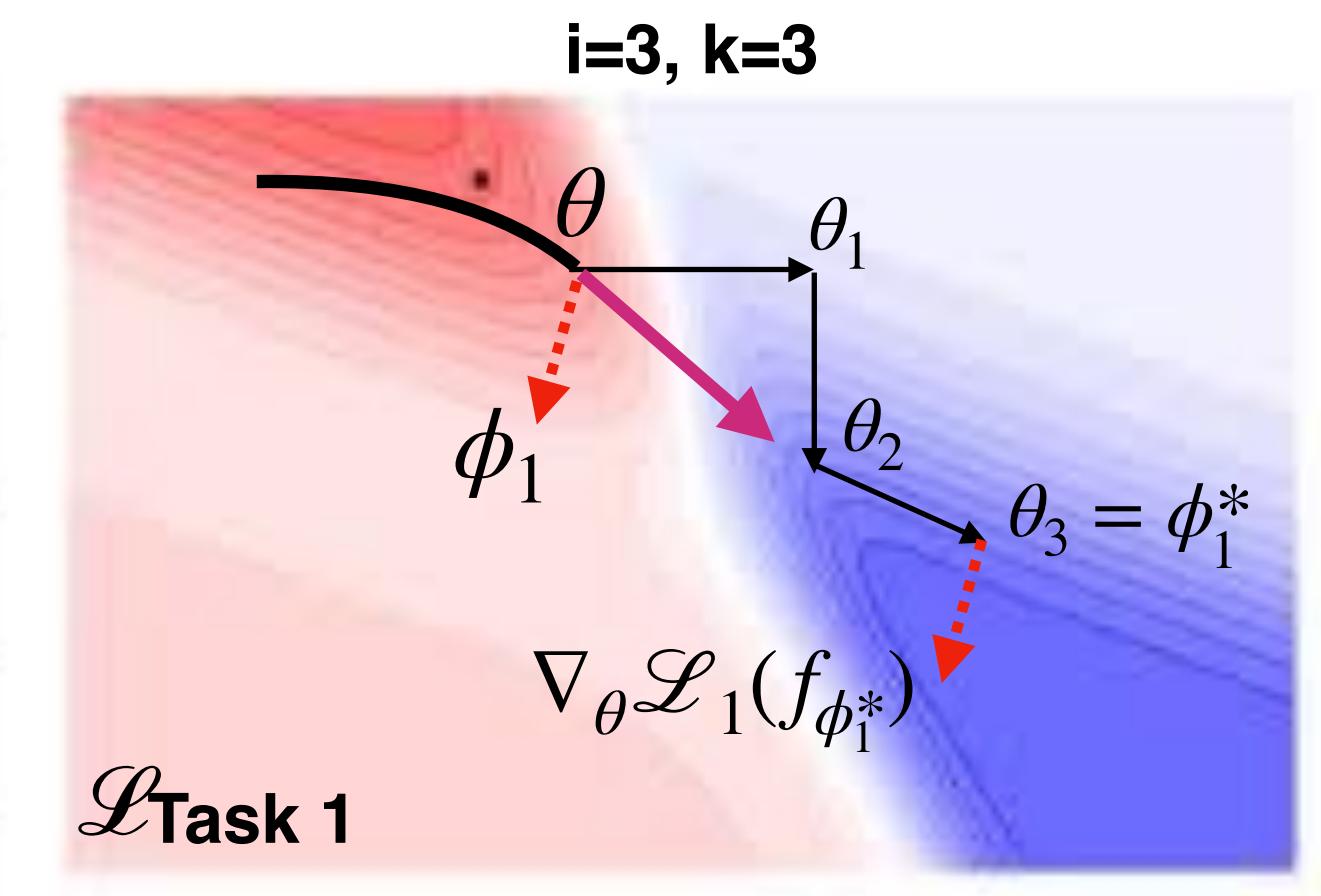
$$\theta \leftarrow \theta - \beta \nabla_\theta \sum_i \mathcal{L}_i(f(\theta - \alpha \nabla_\theta \mathcal{L}_i(f_{\phi_i^*})))$$

meta-gradient: second-order gradient + backpropagate  
compute how changes in  $\Theta$  affect the gradient at new  $\Theta$

## Meta-testing

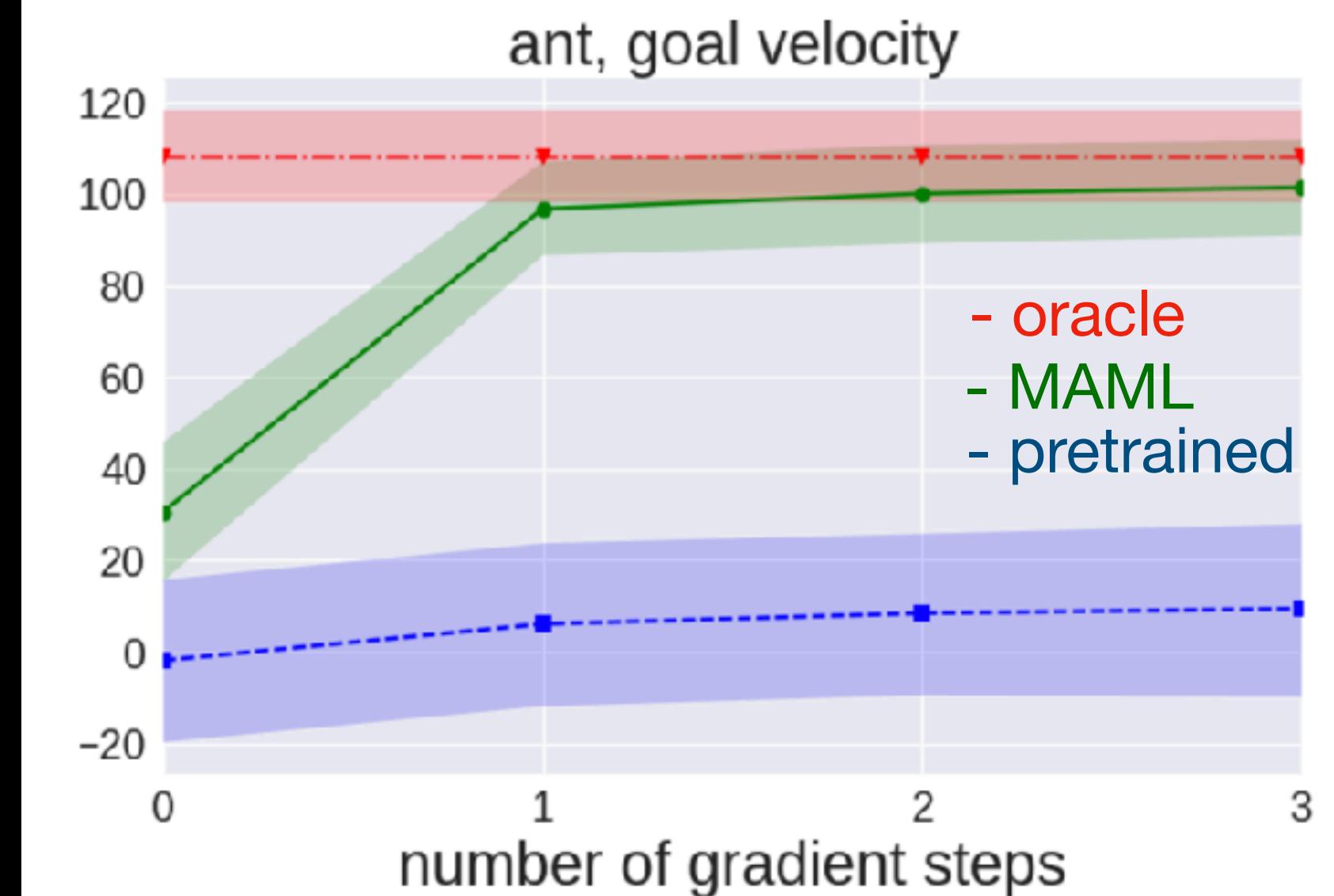
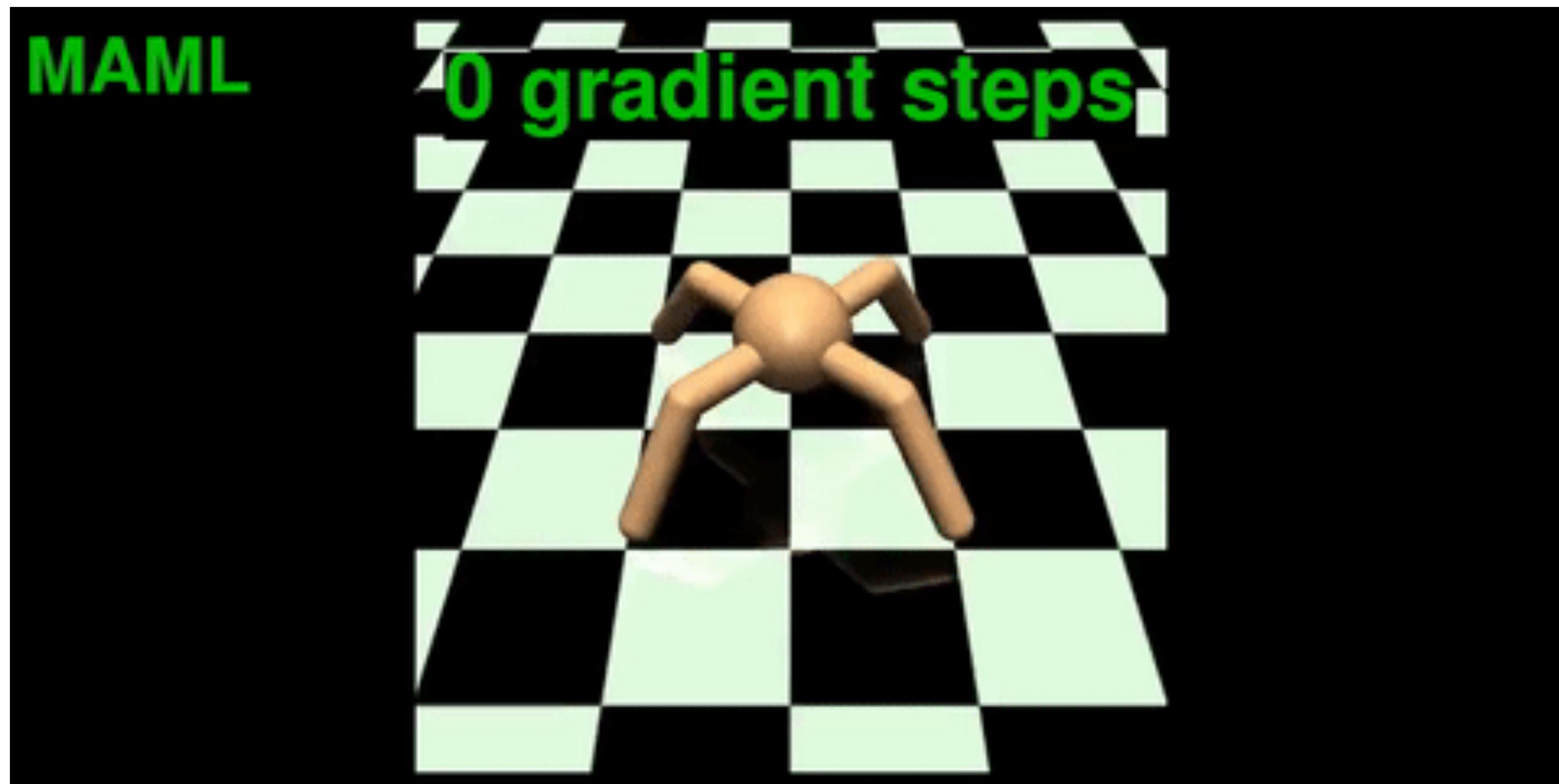
- Training data of new task  $D_{\text{train}}$
- $\Theta^*$ : pre-trained parameters
- Finetune:  $\phi = \theta^* - \alpha \nabla_\theta \mathcal{L}(f_\theta)$

derivative of test-set loss



# Model agnostic meta-learning (MAML)

- Example for reinforcement learning:
  - Goal: reach certain velocity in certain direction



# Other gradient-based techniques

- Changing update rule yield different variants:

- MAML <sup>1,6</sup>

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(f_{\phi_i})$$

$$\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

- MetaSGD <sup>2</sup>

$$\phi_i = \theta - \alpha \text{diag}(w) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

- Tnet <sup>3</sup>

$$\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*}, w)$$

- Meta curvature <sup>4</sup>

$$\phi_i = \theta - \alpha B(\theta, w) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

- WarpGrad <sup>5</sup>

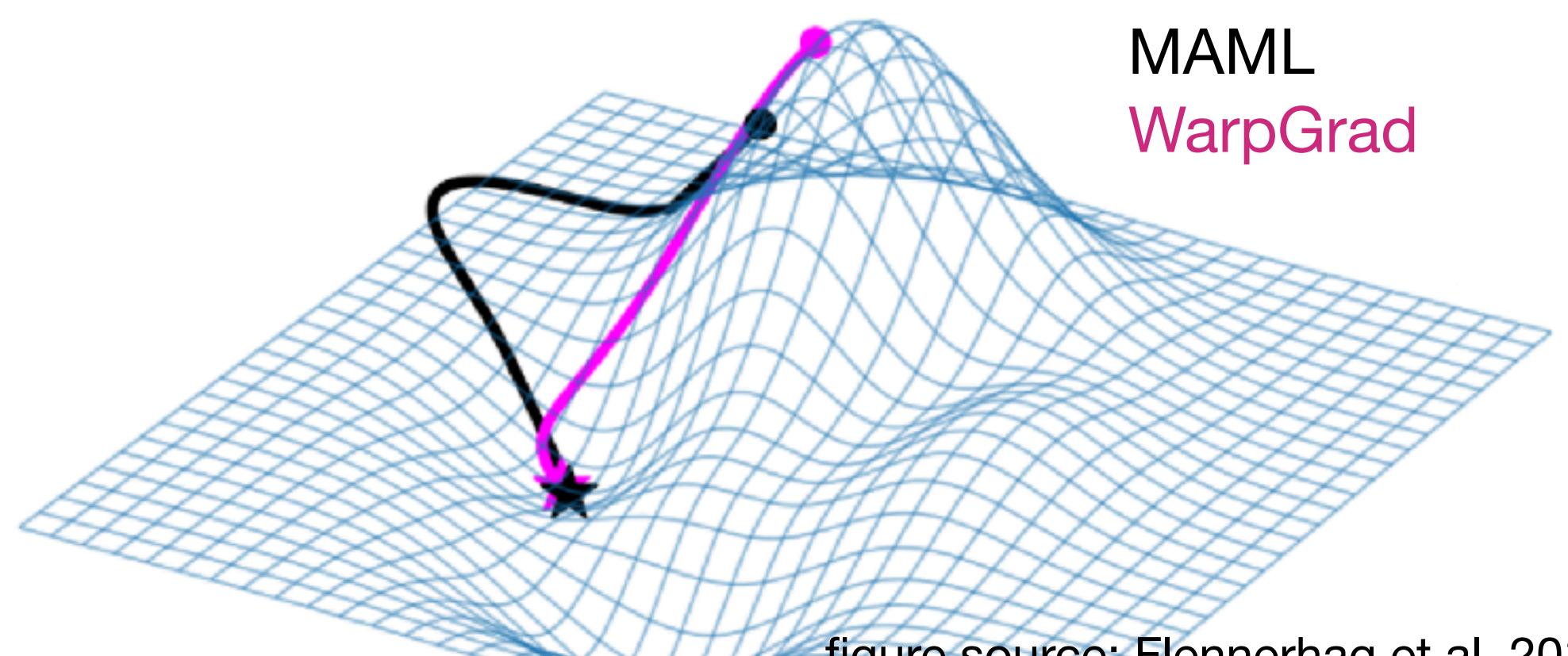
$$\phi_i = \theta - \alpha P(\theta, \phi_i) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

- Online MAML (Follow the Meta Leader) <sup>7</sup>

- Minimizes regret
- Robust, but computation costs grow over time

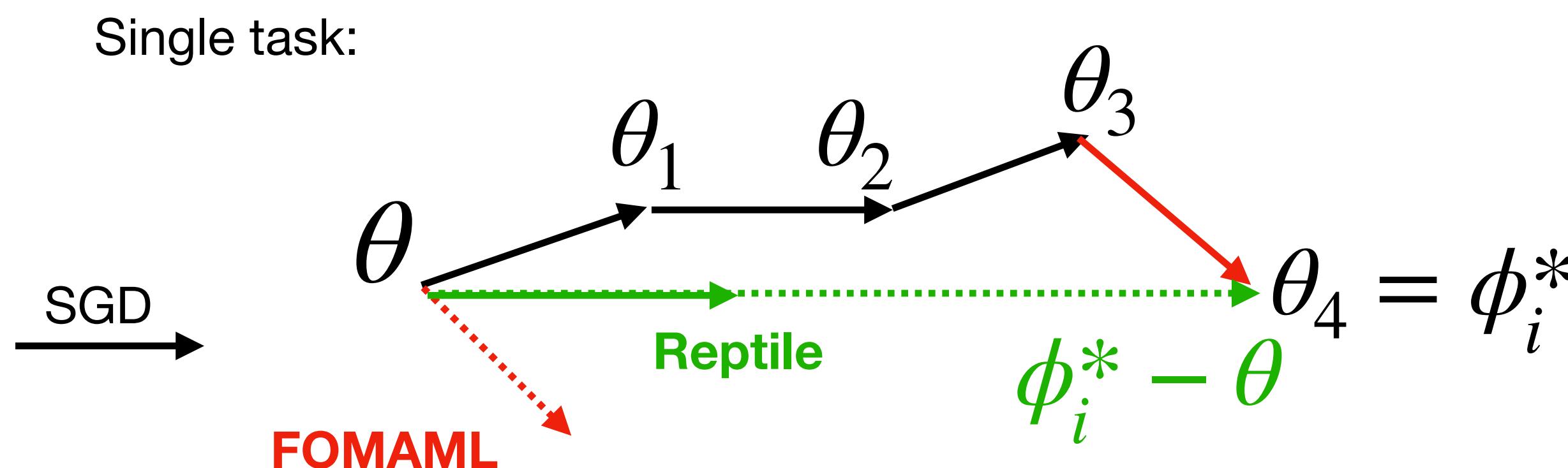
*w: weight per parameter*

**All use second-order gradients,  
Meta-learn a transformation of the  
gradient for better adaptation**

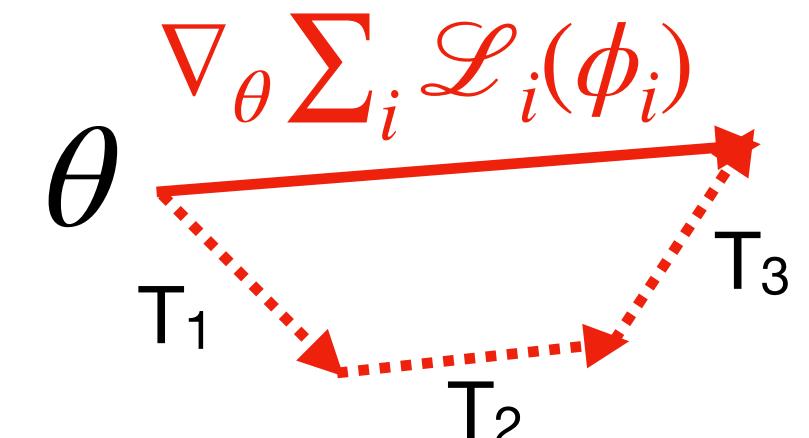


# Scalability

- Backpropagating derivatives of  $\phi_i$  wrt  $\Theta$  is compute + memory intensive (for large models)
- First order approximations of MAML:
  - First order MAML<sup>1</sup> uses only the last inner gradient update:  $\theta \leftarrow \theta - \beta \sum_i \mathcal{L}_i(f_{\phi_i^*})$
  - Reptile<sup>2</sup>: iterate over tasks, update  $\Theta$  in direction of  $\phi_i^*$  :  $\theta \leftarrow \theta - \beta (\phi_i^* - \theta)$

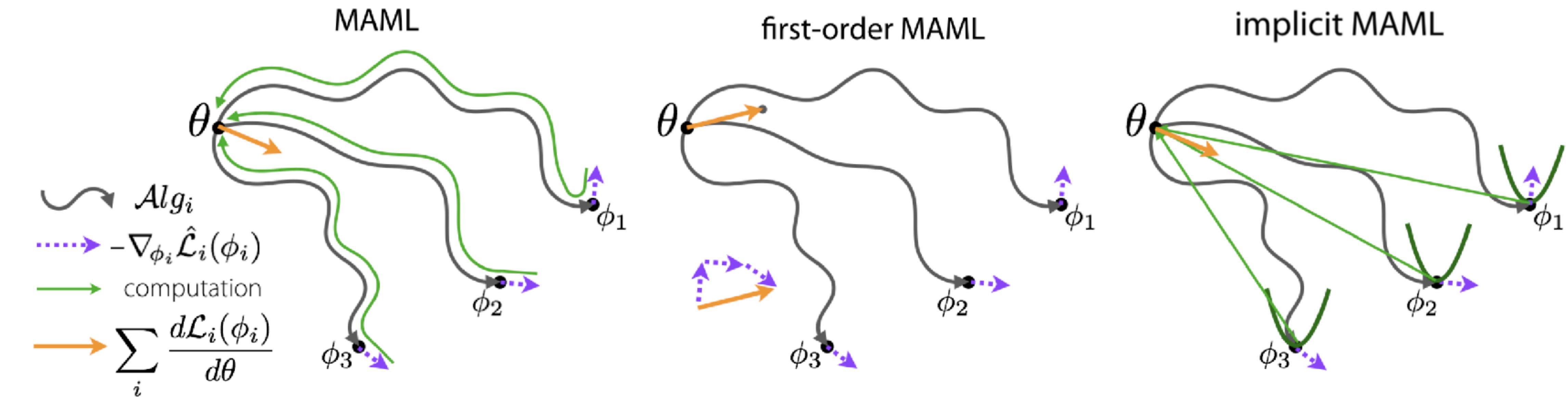


FOMAML meta-gradient:



# Scalability (2)

- Implicit MAML<sup>1</sup>: uses an approximate gradient
  - Compute derivative of  $\phi_i^*$  wrt  $\theta$
  - $\phi_i^*$  could be anywhere. Hence, add penalty:  $\|\phi_i^* - \theta\|^2$
  - Accurate if we stay close to  $\theta$
  - $\mathcal{L}_i(\phi_i) + \lambda \|\phi_i^* - \theta\|^2$  has closed form solution



# Generalizability

<sup>1</sup> Finn et al. 2018

<sup>2</sup> Raghu et al. 2020

<sup>3</sup> Tian et al. 2020

<sup>4</sup> Stadie et al. 2019

- MAML is more resilient to overfitting than many other meta-learning techniques <sup>1</sup>

- Effectiveness seems mainly due to feature reuse (finding  $\Theta$ ) <sup>2</sup>

- Fine-tuning only the last layer equally good

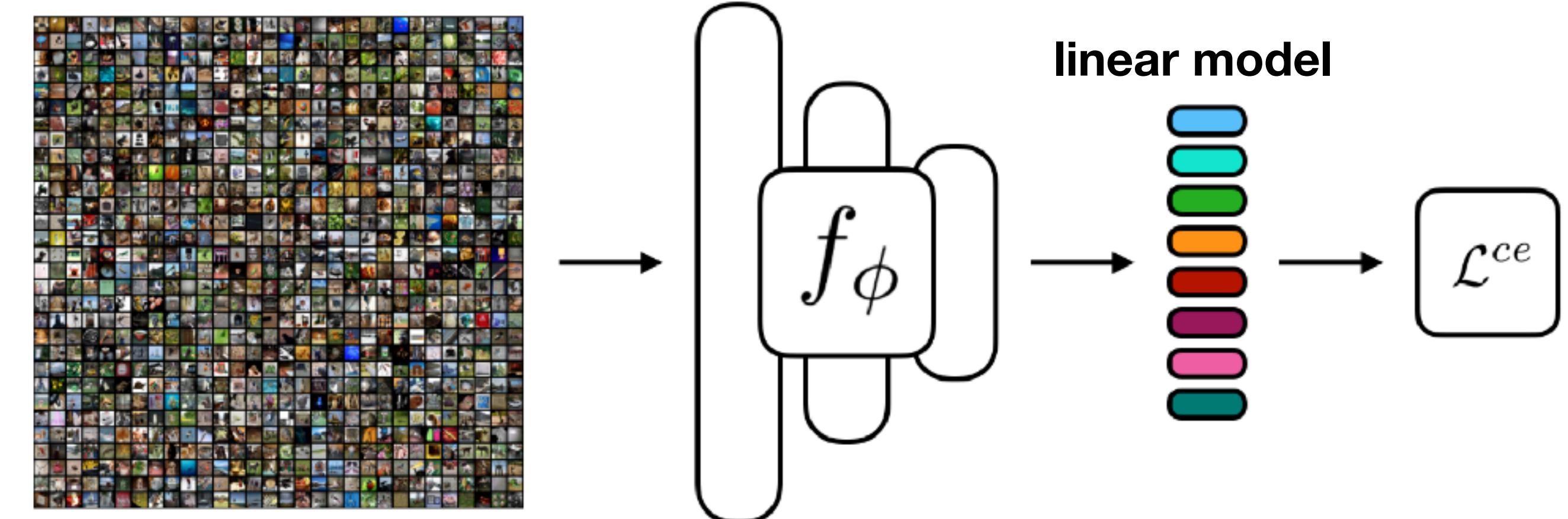
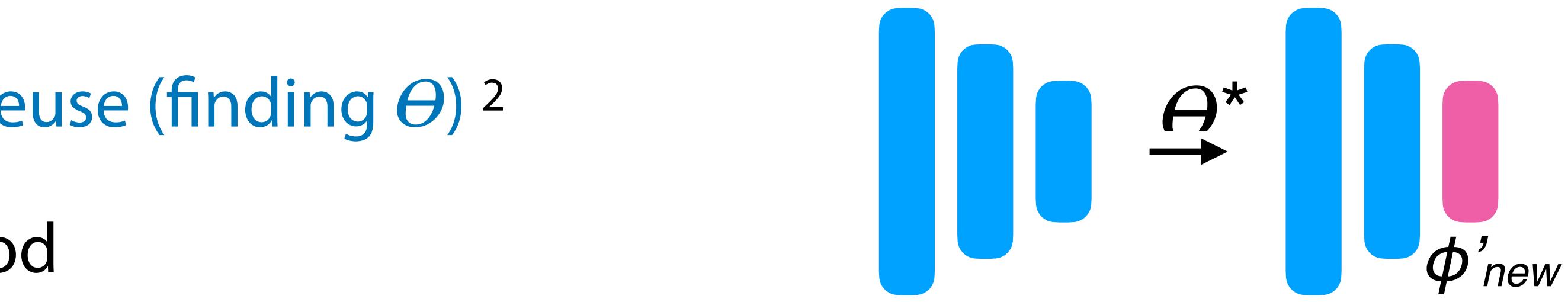
- On few-shot learning benchmarks, a good embedding outperforms most meta-learning <sup>3</sup>

- Learn representation on entire meta-dataset (merged into single task)

- Train a linear classifier on embedded few-shot  $D_{\text{train}}$ , predict  $D_{\text{test}}$

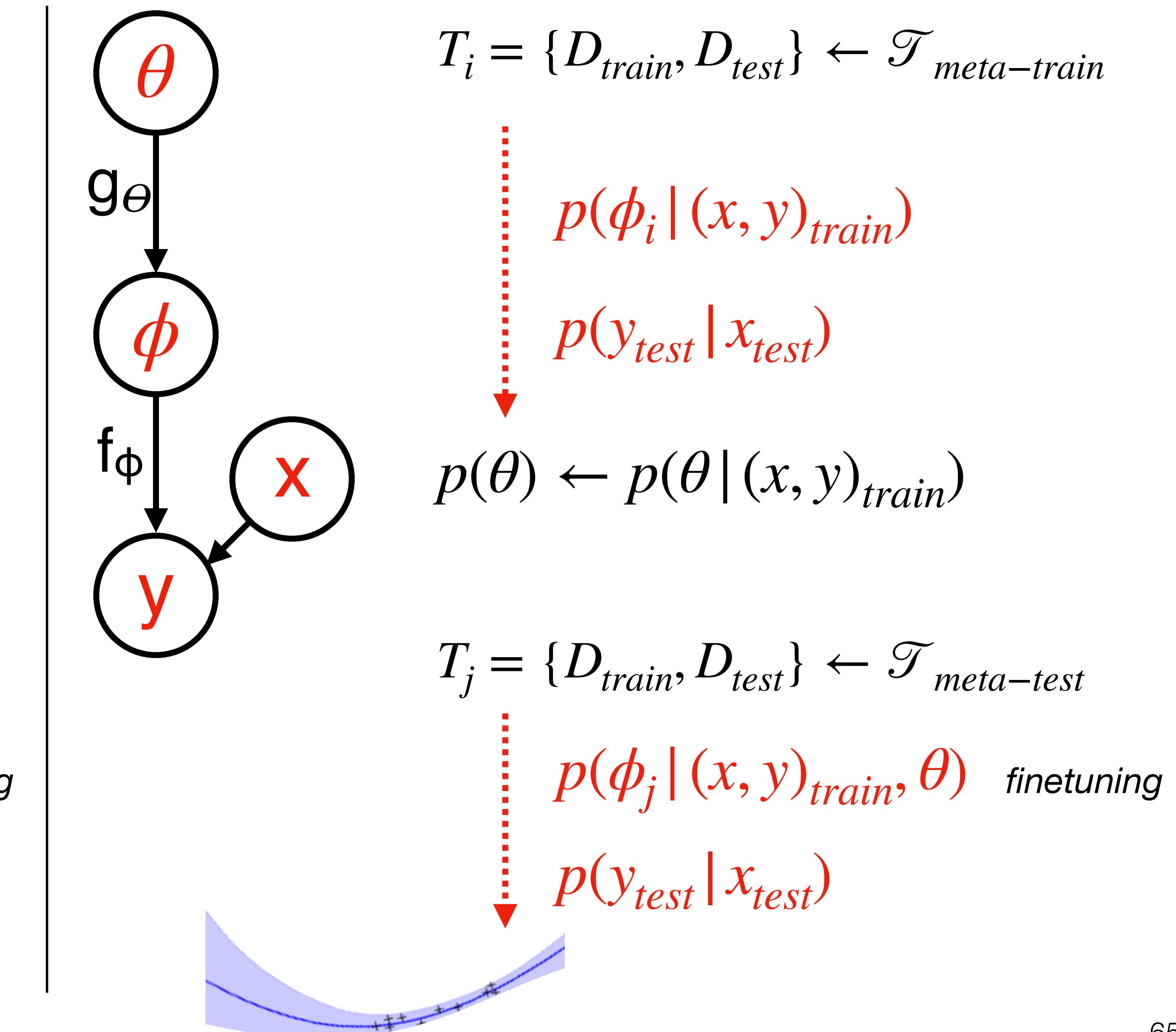
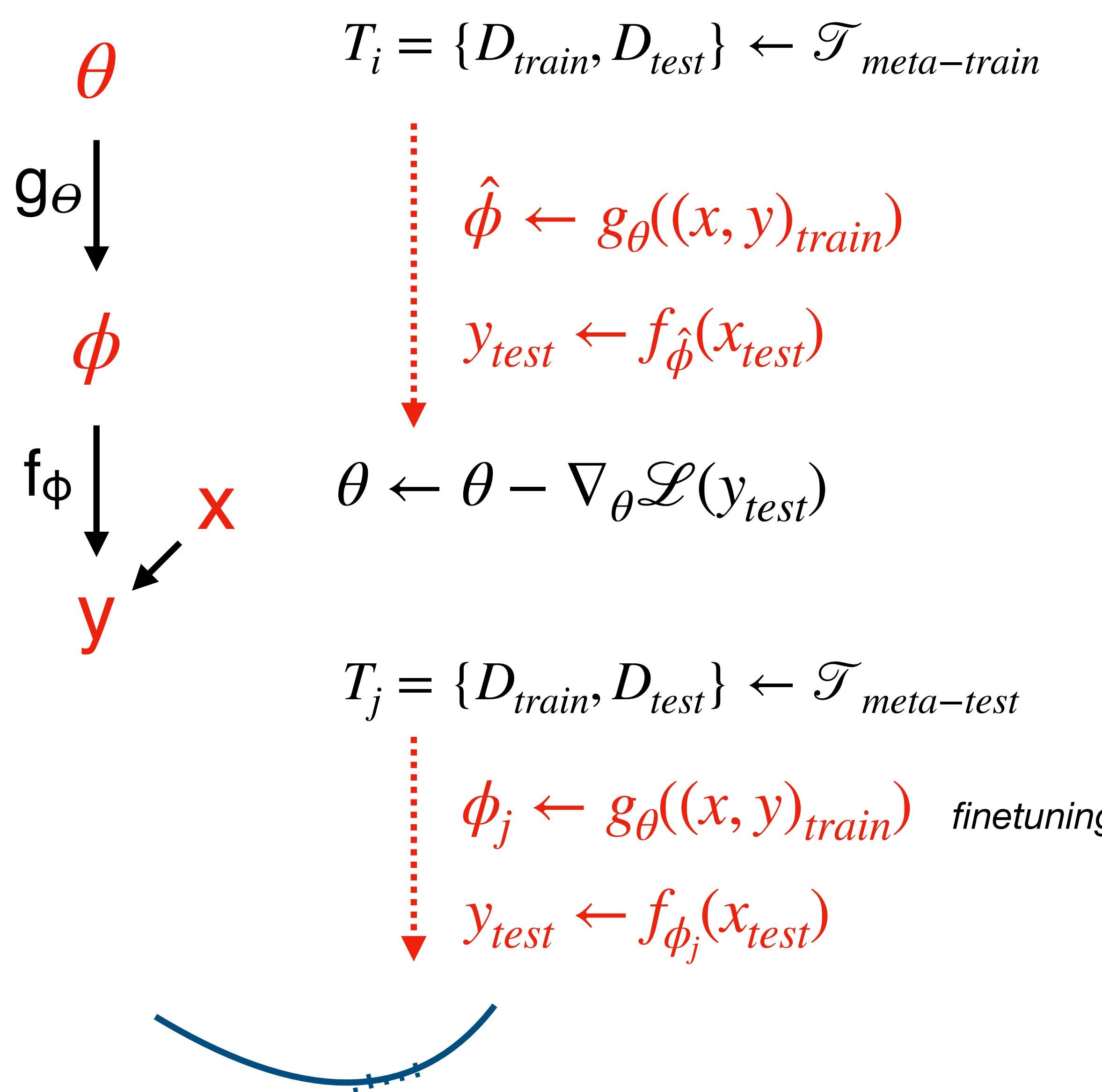
- For meta-RL, also learn how to explore (how to sample new environments)

- E-MAML: add exploration to meta-objective (allows longer term goals) <sup>4</sup>

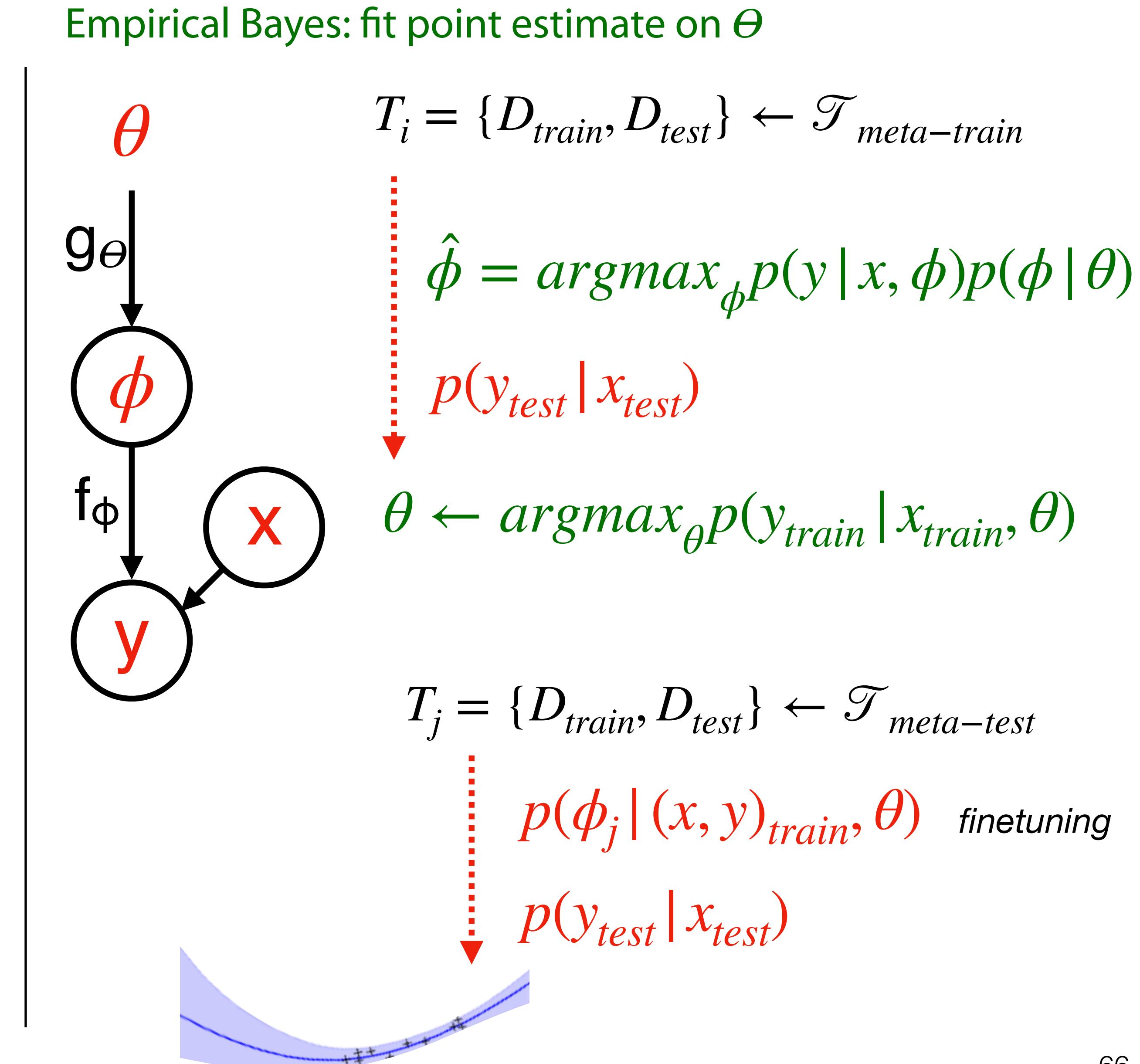
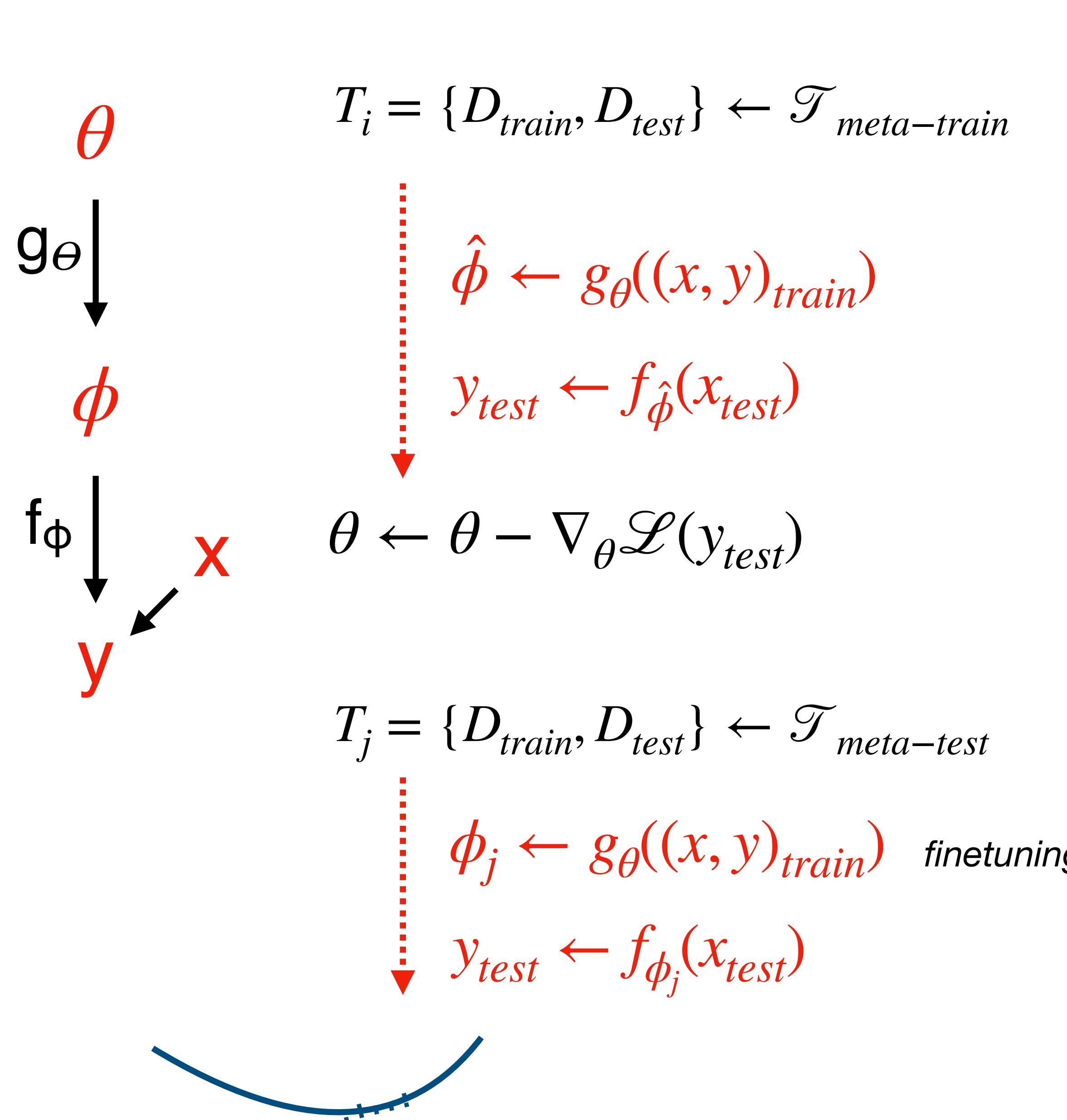


# Bayesian meta-learning

Can meta-learning reason about **uncertainty** in the task distribution?



# Bayesian meta-learning

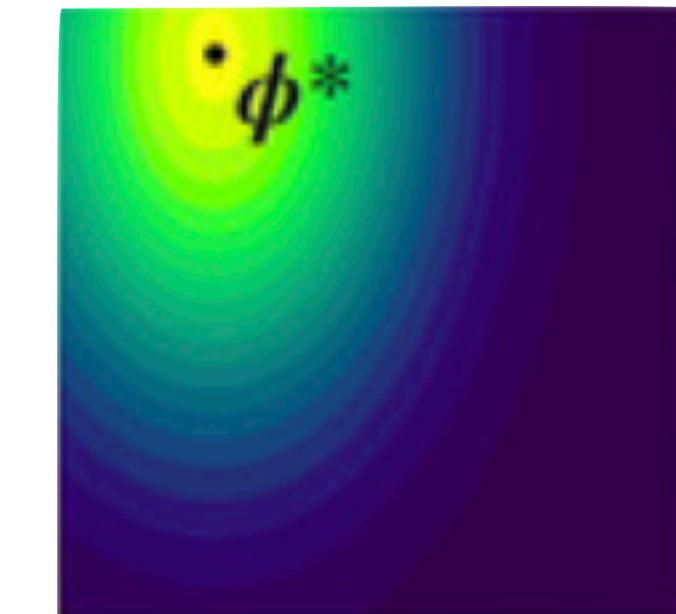


# Bayesian meta-learning

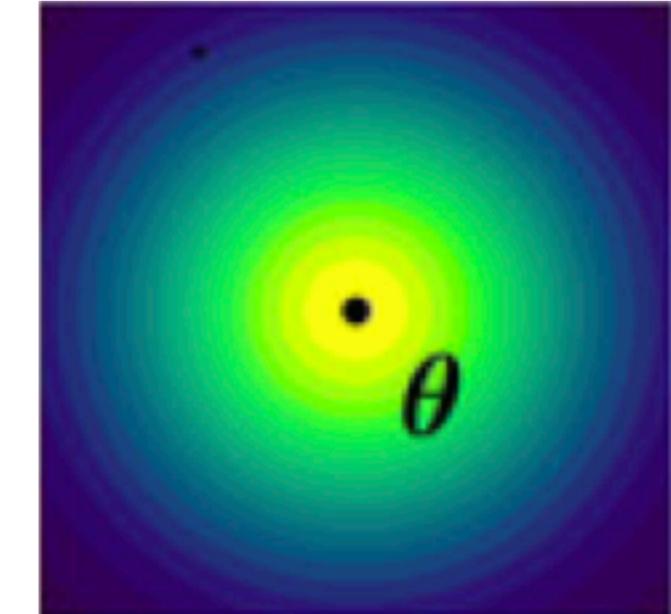
$\phi$  can be factored as the product of a likelihood and a prior

$$\hat{\phi} = \operatorname{argmax}_{\phi} p(y | x, \phi) p(\phi | \theta)$$

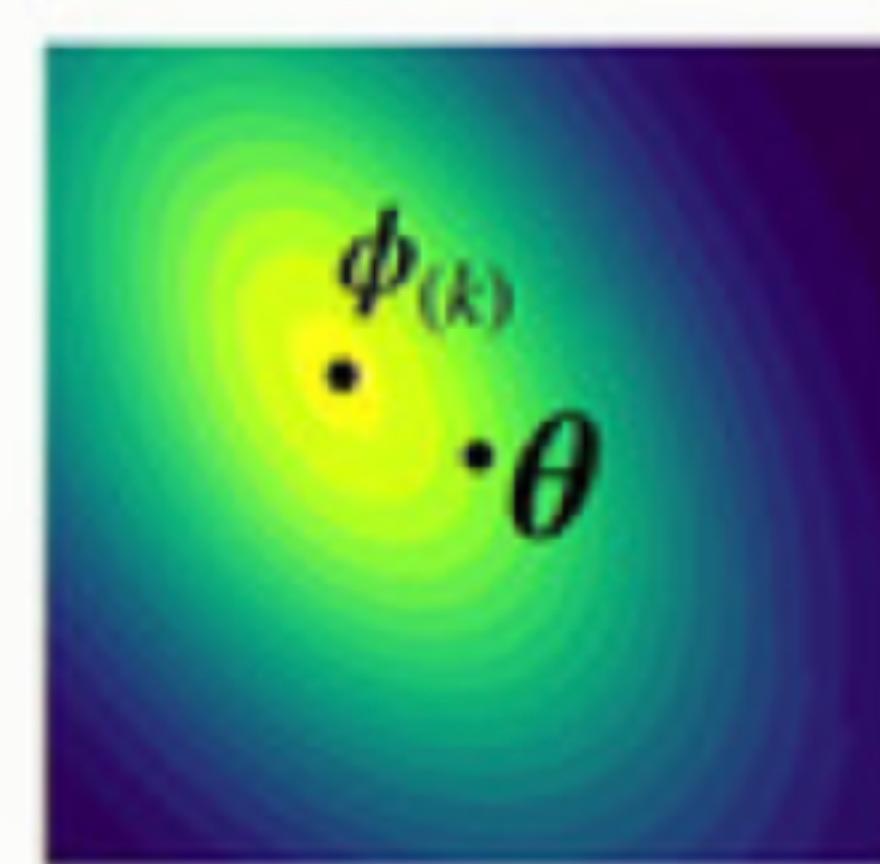
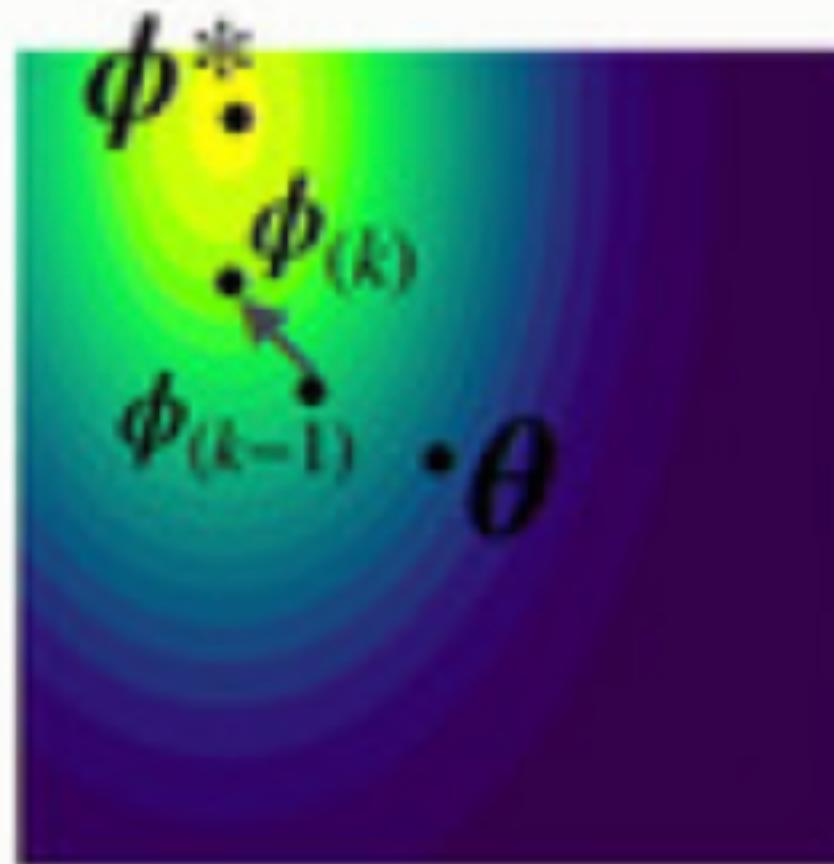
$$p(y | x, \phi)$$



$$p(\phi | \theta)$$



When using  $\Theta$  as the initialization of gradient descent, early stopping is approximately equivalent to placing a prior over the task specific parameters <sup>1</sup>



$\sigma$  depends on k and learning rate

MAML is doing the same

- estimates  $\phi_i^*$  with k steps of gradient descent
- update prior  $\Theta$  to makes it easier to optimize  $\phi_i$

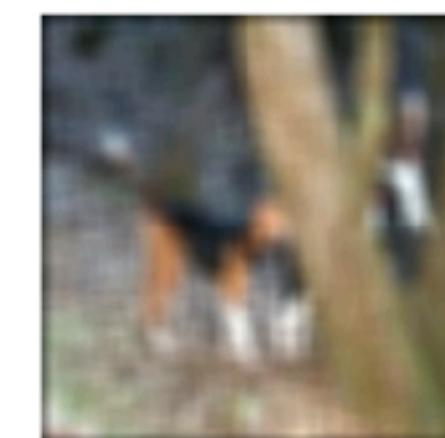
# Fully Bayesian meta-learning

- Alternatively, use approximation methods to represent uncertainty over  $\phi_i$
- Sampling technique + variational inference: PLATIPUS<sup>1</sup>, BMAML<sup>2</sup>, ABML<sup>3</sup>
- Laplace approximation: LLAMA<sup>4</sup>
- Variational approximation of posterior:
  - Neural Statistician<sup>5</sup>, Neural Processes<sup>6</sup>
- What if our tasks are not IID?
  - Impose additional structure, e.g. with task-specific variable z<sup>7</sup>

mini-ImageNet with filters for non-homogeneous tasks:



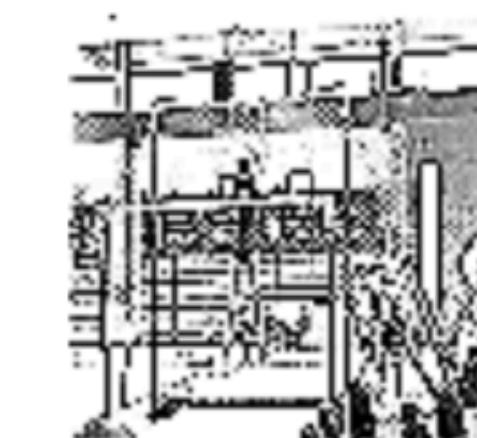
(a) plain



(b) blur



(c) night



(d) pencil

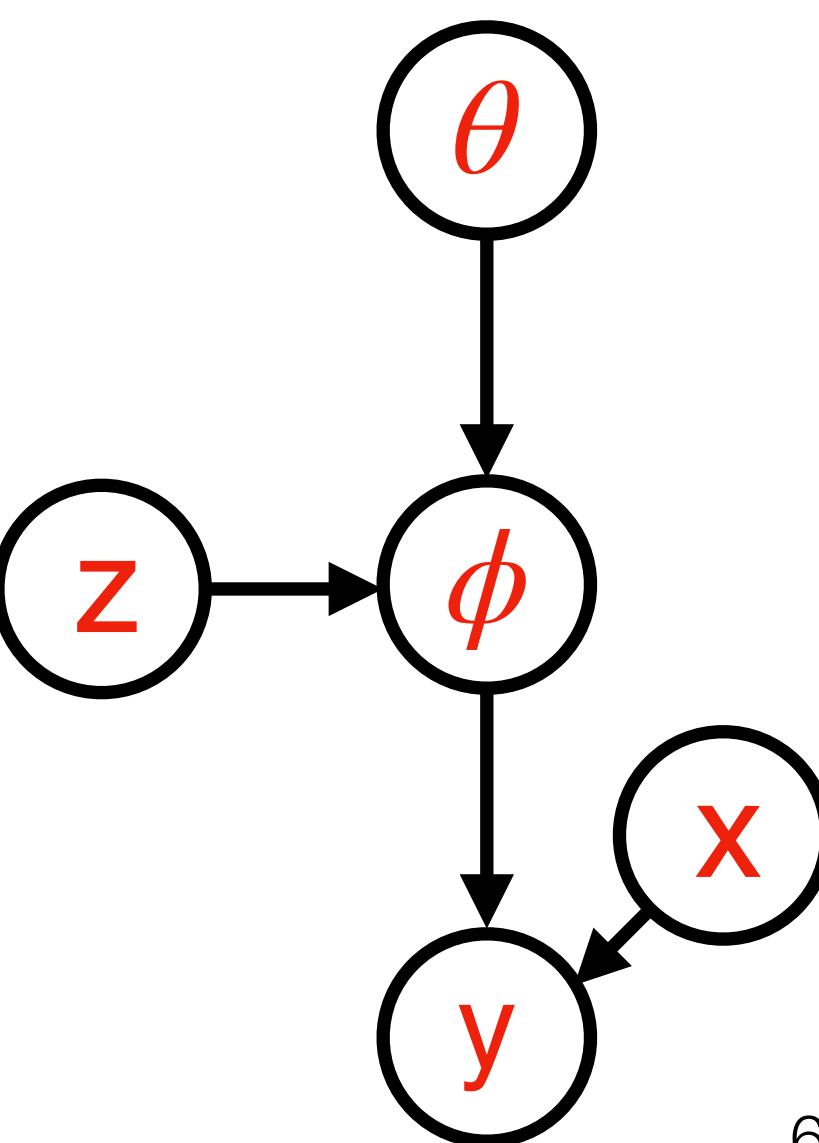
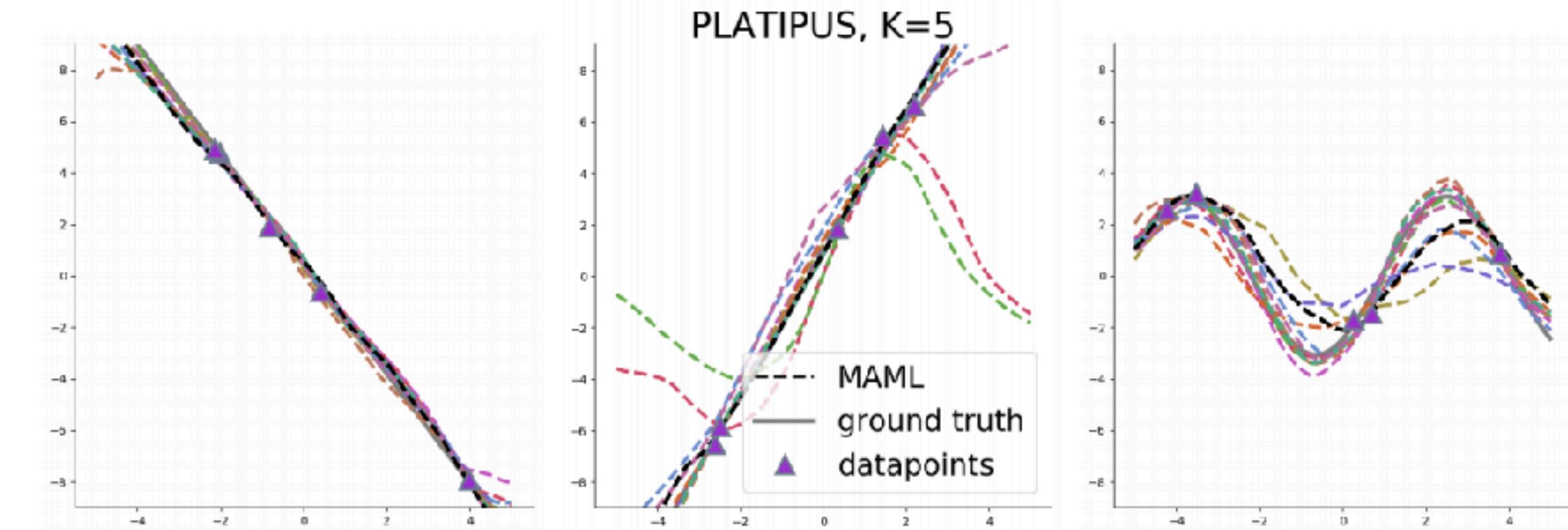


Figure source: [Jerfel et al. 2019](#)

<sup>1</sup> [Finn et al. 2019](#)

<sup>2</sup> [Kim et al. 2018](#)

<sup>3</sup> [Ravi et al. 2019](#)

<sup>4</sup> [Grant et al. 2018](#)

<sup>5</sup> [Edwards et al. 2017](#)

<sup>6</sup> [Garnelo et al. 2018](#)

<sup>7</sup> [Jerfel et al. 2019](#)

# Meta-learning optimizers

<sup>1</sup> Bengio et al. 1995

<sup>2</sup> Schmidhuber 1992

<sup>3</sup> Runarsson and Jonsson 2000

<sup>4</sup> Hochreiter 2001

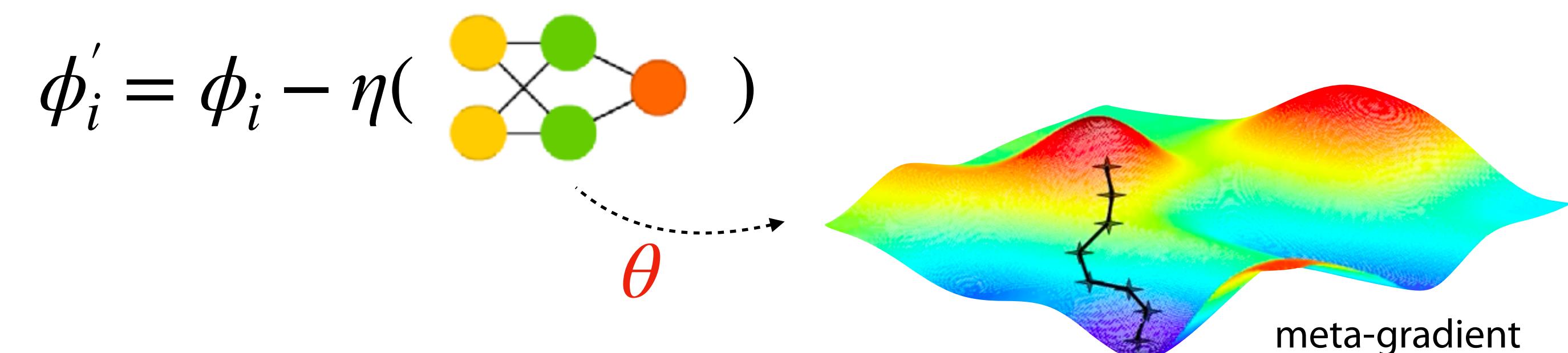
- Our brains probably don't do backprop, instead:

- Simple bio-inspired rules to update weights <sup>1</sup>

$$\phi'_i = \phi_i - \eta y_{pre(i)} k$$

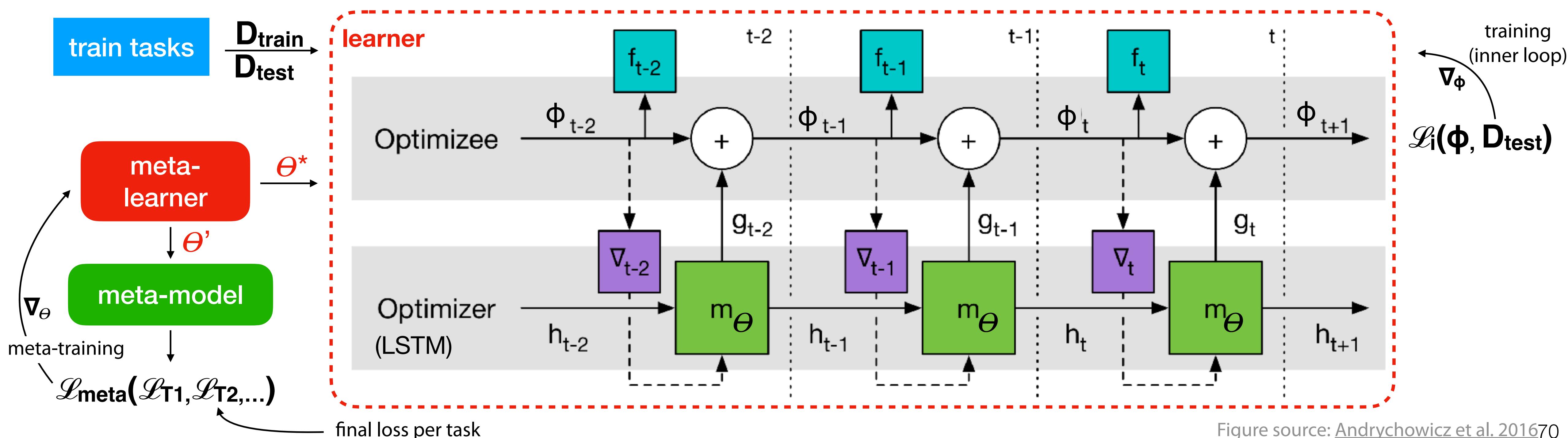
presynaptic activity  
learning rate      reinforcing signal

- Fast weights <sup>2</sup>: networks that continuously modify the weights of another network
  - Gradient based <sup>3, 4</sup>: parameterize the update rule using a neural network
    - Learn meta-parameters across tasks, by gradient descent



# Meta-learning optimizers

- Represent update rule as an LSTM <sup>1,2,3</sup>, hierarchical RNN <sup>4</sup>
  - Optimizee: receives weight update  $g_t$  from meta-learned optimizer
  - Optimizer: receives gradient estimate  $\nabla_t$  from optimizee
- Meta-learner learns optimizer parameters with gradient descent across tasks
  - e.g. Image classification <sup>2,4</sup>, few-shot learning <sup>3</sup>



# Meta-learning optimizers

- Meta-learned (RNN) optimizers ‘rediscover’ momentum, gradient clipping, learning rate schedules, learning rate adaptation, ... <sup>1</sup>

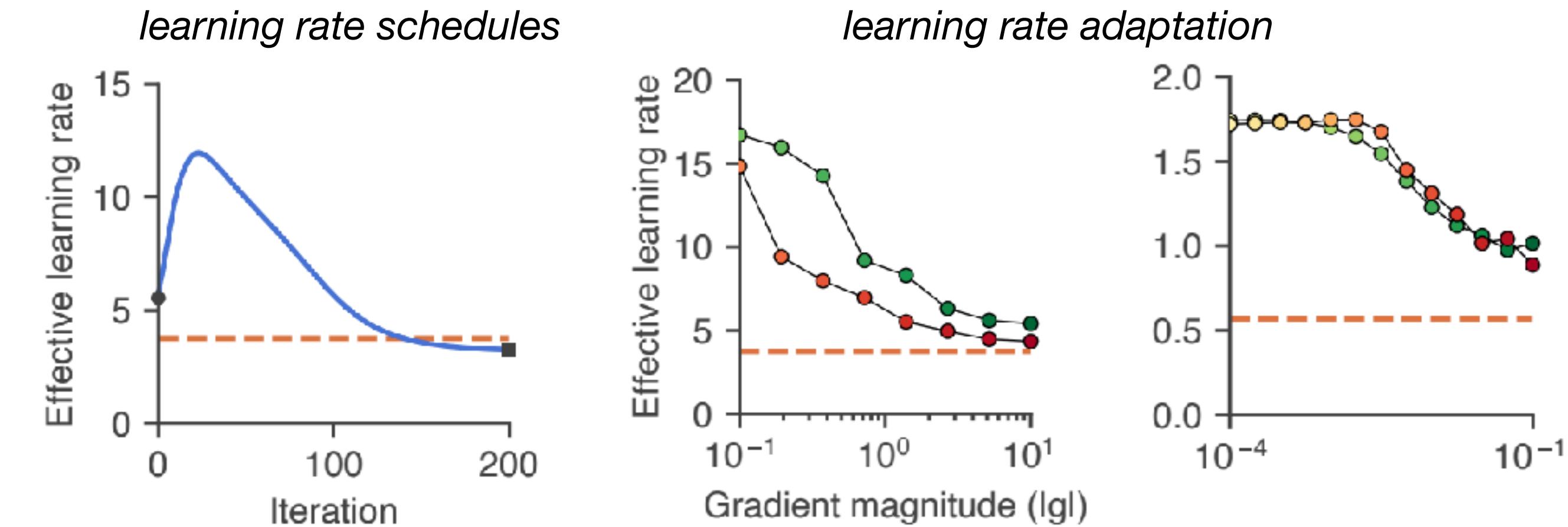
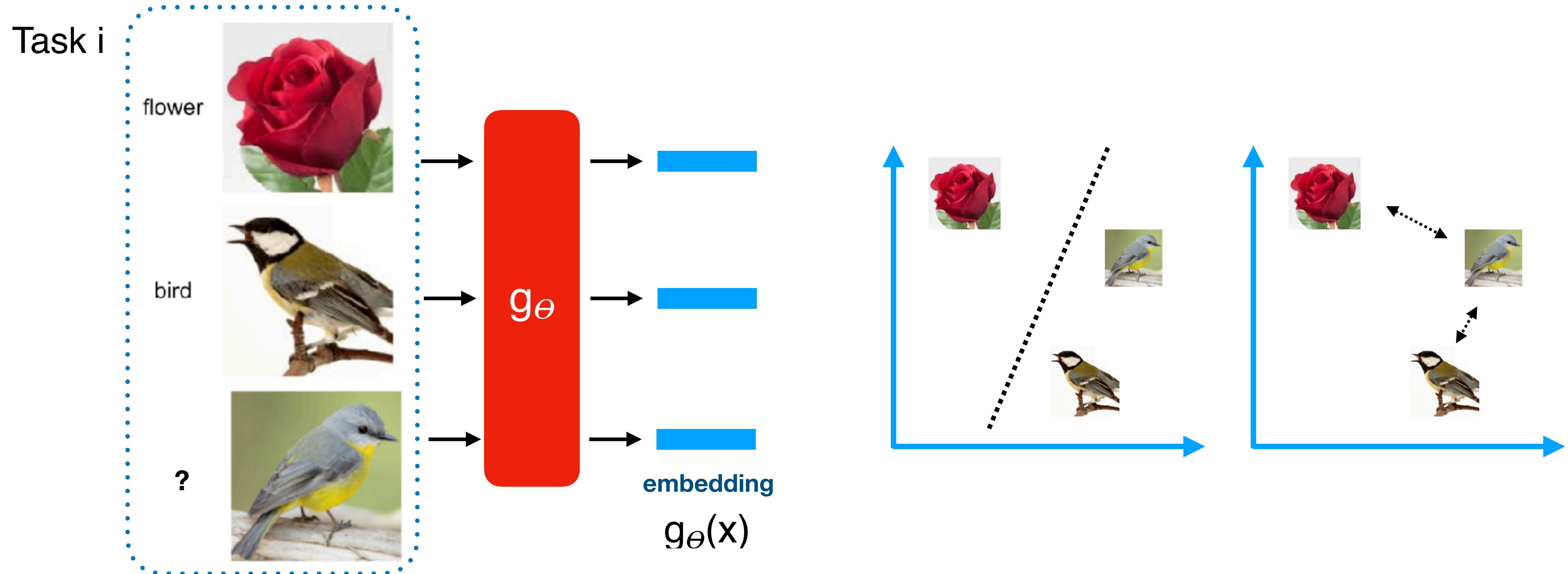


Figure source:  
Maheswaranathan et al. 2020

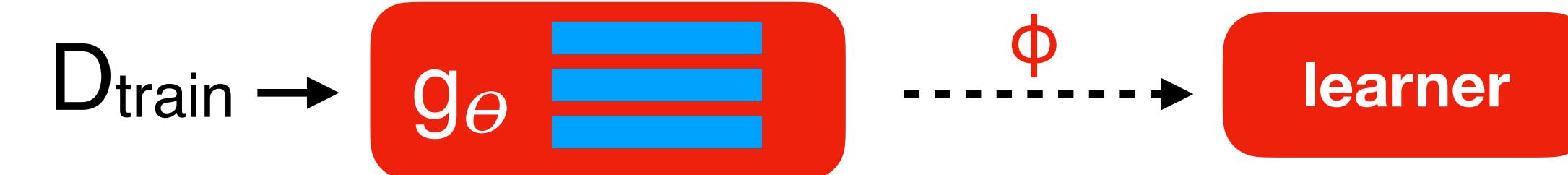
- RL-based optimizers: represent updates as a policy, learn using guided policy search <sup>2</sup>
- Combined with MAML:
  - learn per-parameter learning rates <sup>3,4</sup>
  - learn precondition matrix (to ‘warp’ the loss surface) <sup>5,6</sup>
- Black-box optimizers: meta-learned with an RNN <sup>7</sup>, or with user-defined priors <sup>8</sup>
- Speed up backpropagation by meta-learning sparsity and weight sharing <sup>9</sup>

# Metric learning

Learn an embedding network  $\Theta$  that transforms data  $\{D_{train}, D_{test}\}$  across all tasks to a representation that allows easy similarity comparison



Can be seen as a simple black box meta-model  
• Often non-parametric (independent of  $\phi$ )



# Matching networks

- Classifier: probability distribution based on similarity between  $x_{test}$  and  $x_i \in D_{train}$

$$p(y | x_{test}, D_{train}) = \sum_{i=1}^k a(x_{test}, x_i) y_i$$

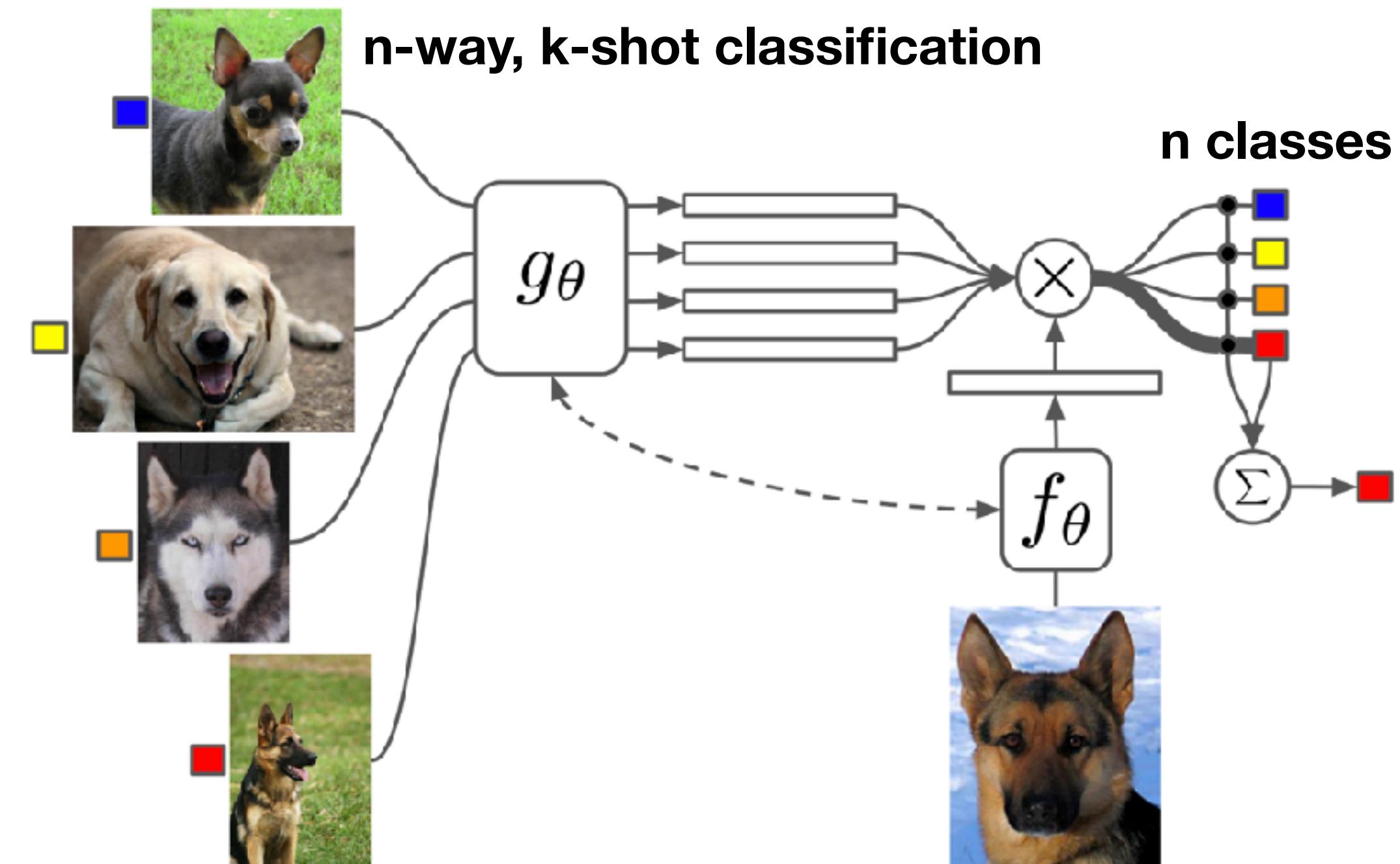
- Similarity: attention kernel based on cosine distance between embedded data points

$$a(x, x_i) = \text{softmax}(\cos(f(x), g(x_i)))$$

- Learn two embeddings:  $\theta \leftarrow \theta - \alpha \nabla_\theta \mathcal{L}(y, \hat{y})$

- $g_\theta$  for  $D_{train}$  samples
- $f_\theta'$  for  $D_{test}$  samples

- Two options:
- Simple embedding:  $\Theta=\Theta'$ , so  $g_\theta=f_\theta'$
- Contextual embedding:  $g$  is a bidirectional LSTM and  $f$  is an attention LSTM



# Prototypical networks

- Use an embedding function  $f_\theta$  to encode each data point
- Define a prototype  $v_c$  for every class  $c$  based on the examples of that class  $D_{train}^{(y)}$

$$v_c = \frac{1}{|D_{train}^{(y)}|} \sum_{(x_i, y_i) \in D_{train}^{(y)}} f_\theta(x_i)$$

- Class distribution for input  $x$  is based on inverse distance between  $x$  and prototypes

$$p(y = c | x_{test}) = \text{softmax}(-d_\phi(f_\theta(x), v_c))$$

- Distance function can be any differentiable distance
  - E.g. squared Euclidean
- Loss function to learn the embedding:

$$\mathcal{L}(\theta) = -\log p_\theta(y = c | x)$$

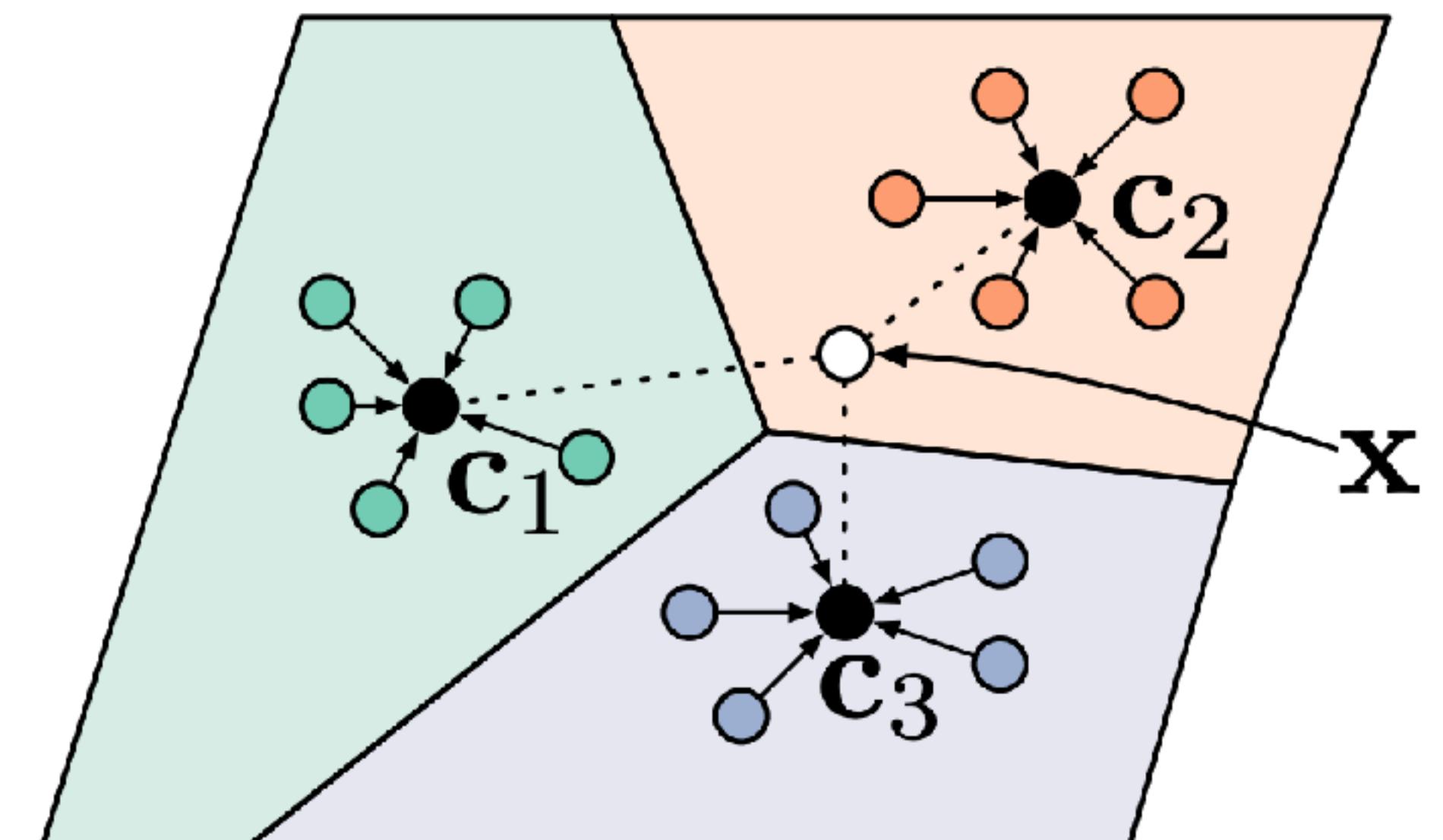
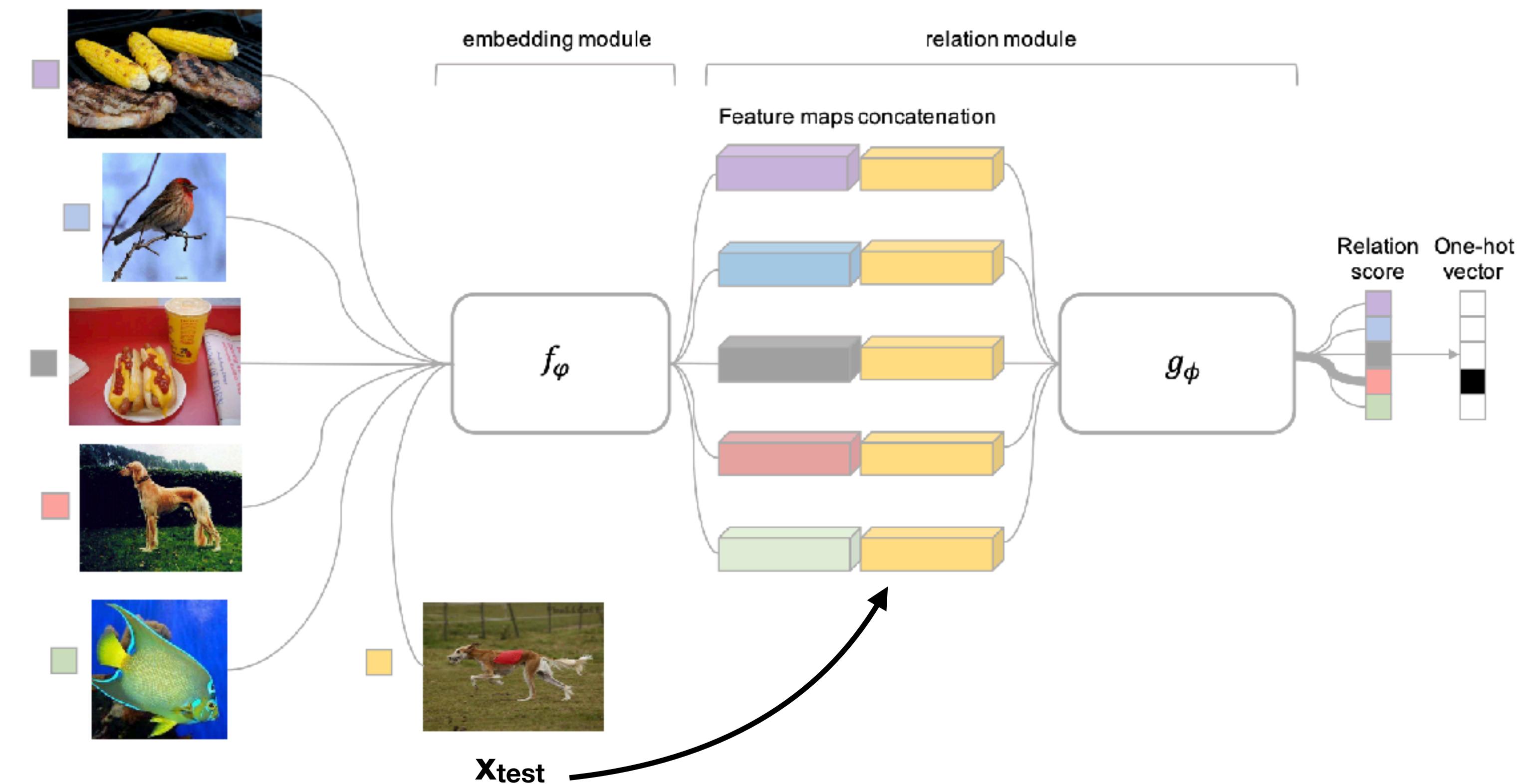


Figure source: Snell et al. 2017 74

# Relation networks

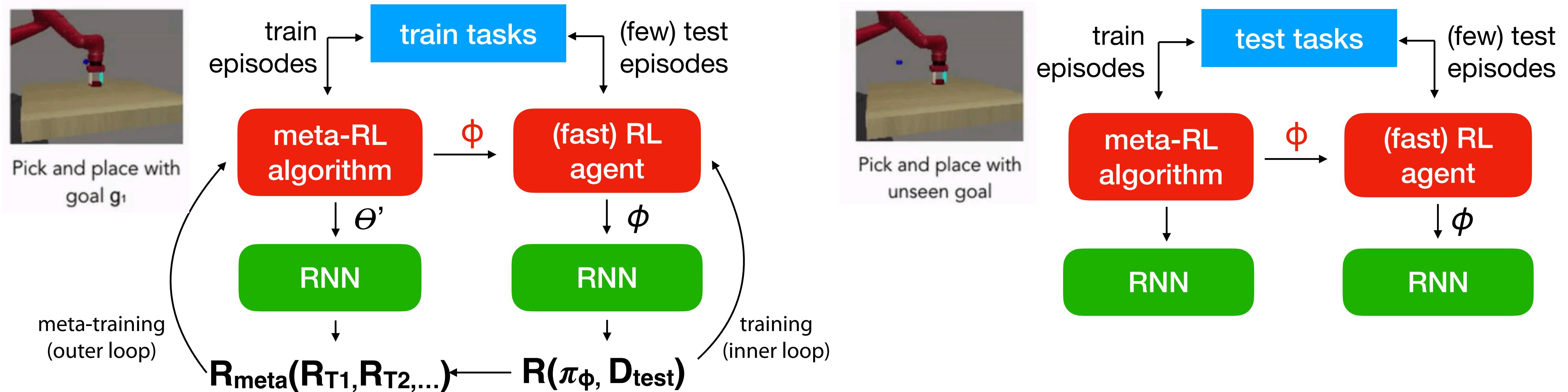
- Similar to matching networks, but with a trainable similarity metric
    - Learns a non-linear relationship between the data points
  - Learn an embedding network  $f_\theta$  and a relation network  $g_\theta'$
  - The relationship/similarity between a pair of inputs:
- $$r_{ij} = g_\theta'(concat(f_\theta(x_i), f_\theta(x_j)))$$
- Predictions based on the examples most related to  $x_{test}$
  - More expressive power, often beats other metric learners



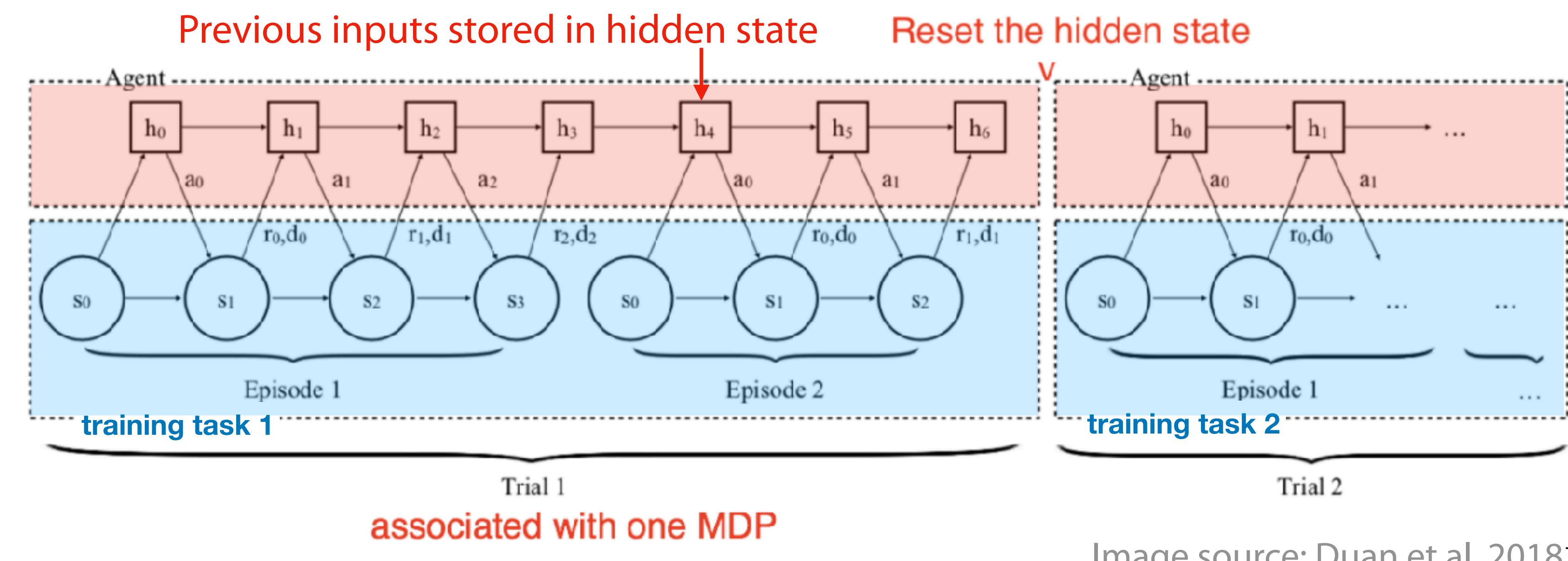
# Metric learning

- Quite a few other techniques exist
  - Siamese neural networks <sup>1</sup>
  - Graph Neural Networks <sup>2</sup>
    - Also applicable for semi-supervised and active learning
  - Attentive Recurrent Comparators <sup>3</sup>
    - Compares inputs not as a whole but by parts (e.g. image patches)
  - MetaOptNet <sup>4</sup>
    - Learns embeddings so that linear models can distinguish between classes
- Overall:
  - Fast at test time, although pair-wise comparisons limit task size
  - Mostly limited to few-shot supervised tasks
  - Fails when test tasks are more distant: no way to adapt

# Black-box model for meta-RL



- RNNs serve as dynamic task embedding storage
- Maximize expected reward in each trial
- Very expressive, perform very well on short tasks
- Longer horizons are an open challenge



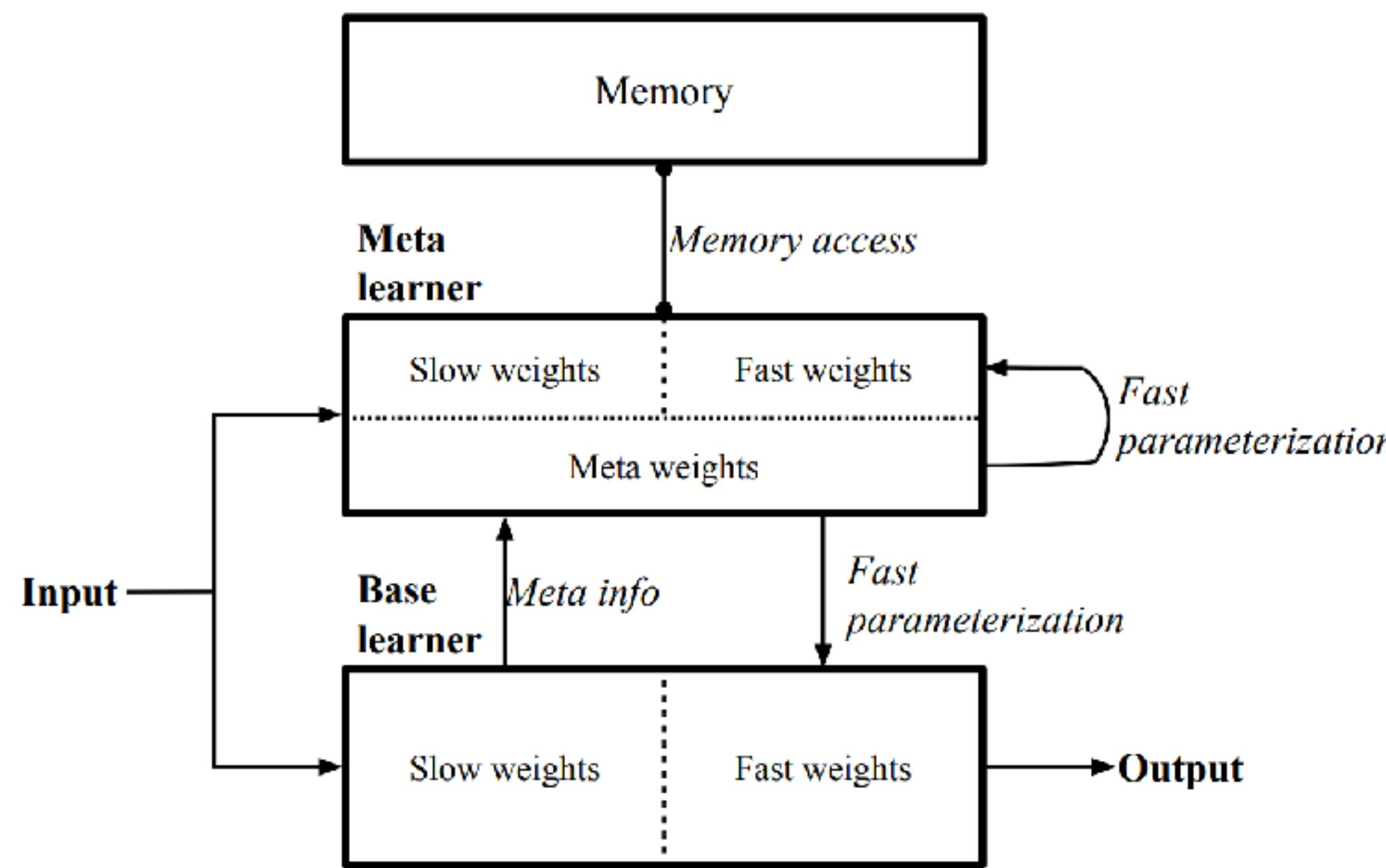
# Other black-box models

[Santoro et al. 2016](#)

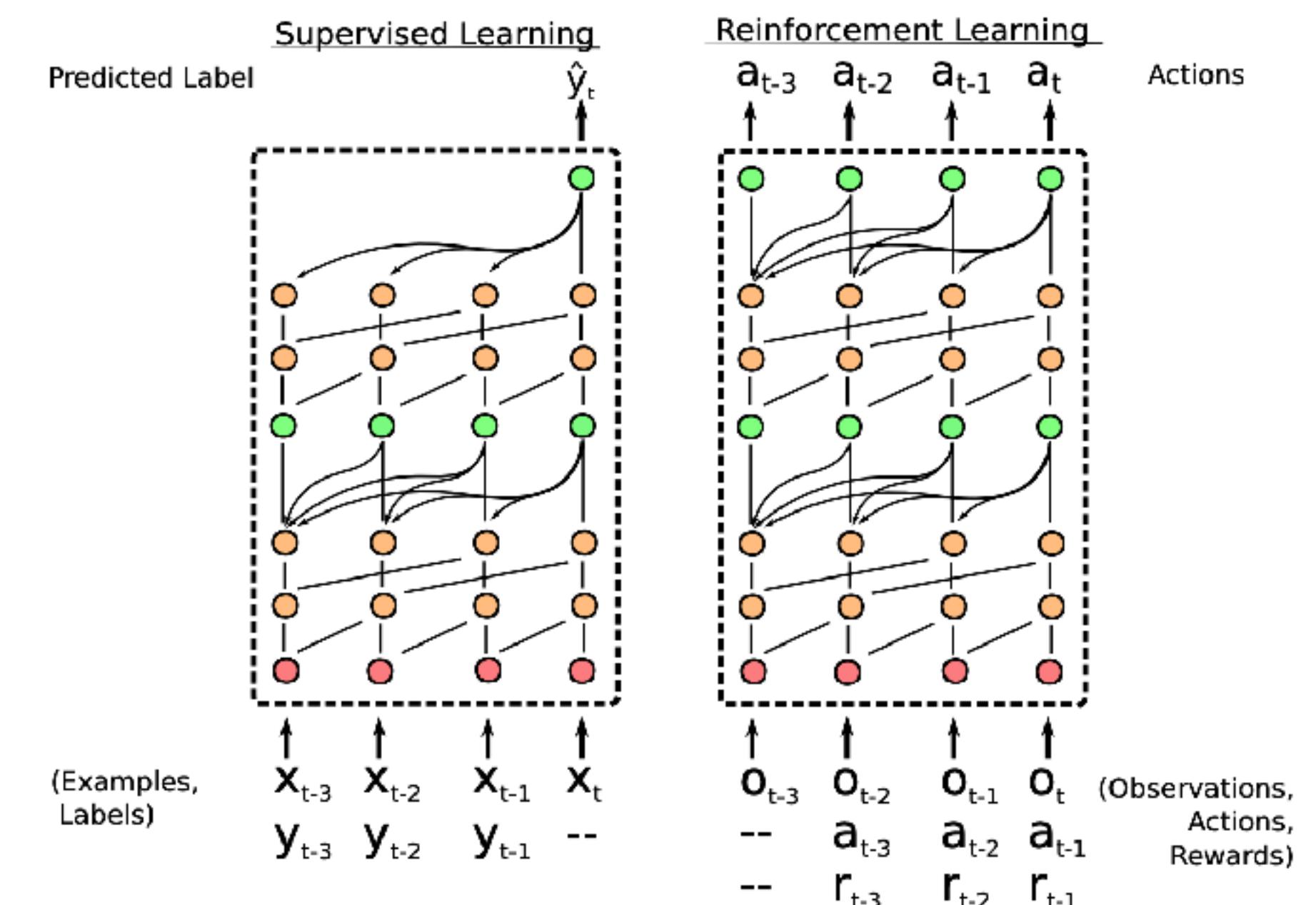
[Munkhdalai et al. 2017](#)

[Mishra et al. 2018](#)

- Memory-augmented NNs <sup>1</sup>
  - Uses neural Turing machines: short term + long term memory
- Meta Networks <sup>2</sup>
  - Meta-learner that returns ‘fast weights’ for itself and the base network solving the task
- Simple Neural attentive meta- learner (SNAIL)
  - Aims to overcome memory limitations of RNNs with series of 1D convolutions



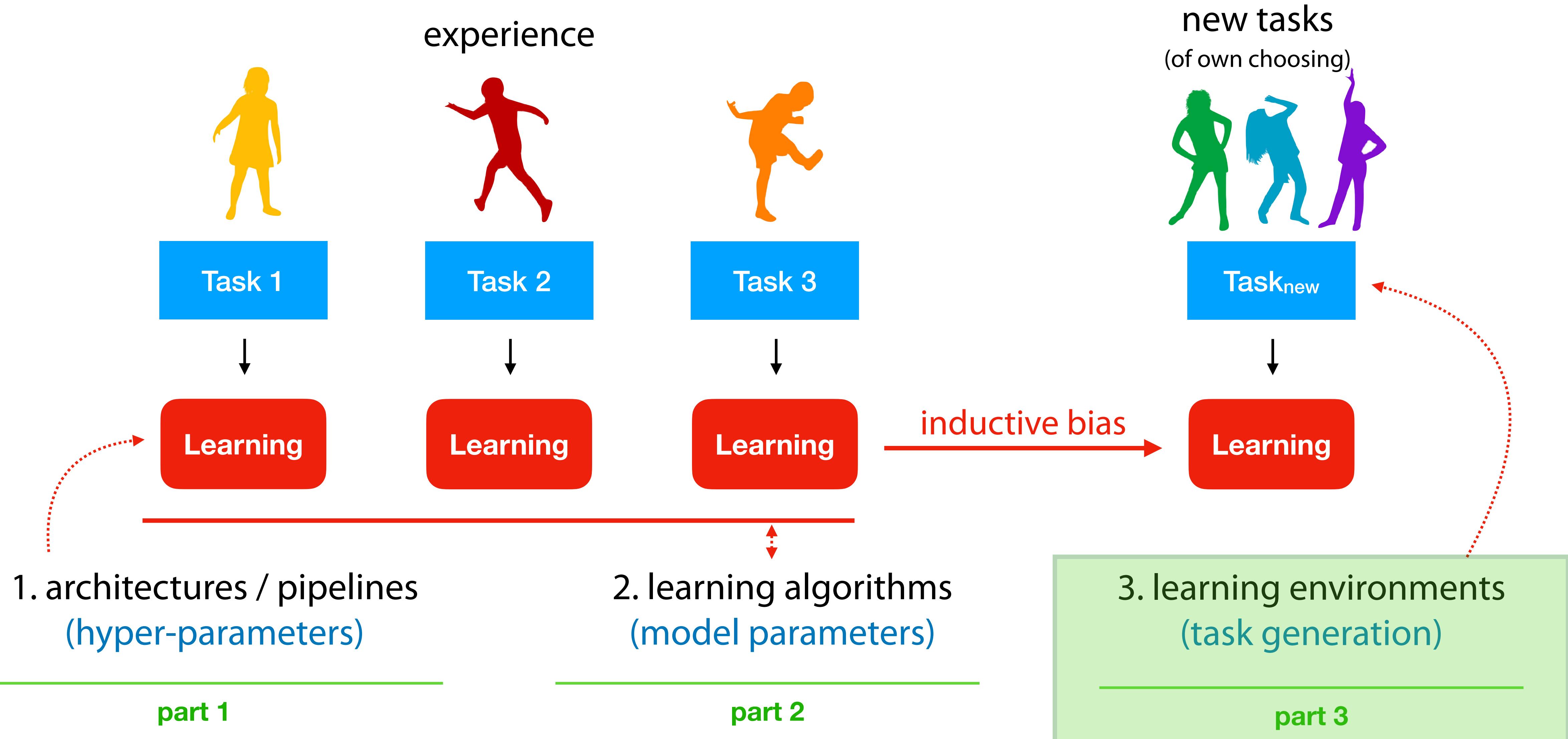
Memory-augmented NNs. Image source: Munkhdalai et al. 2017



SNAIL. Image source: Mishra et al. 2018

# What can we learn to learn?

*3 pillars*



# Training Task Acquisition

- Ultimately, meta-learning translates constraints on the learner to constraints on the data
  - The biases we don't put in manually have to be learnable from data
  - Can we automatically create new tasks to inform and challenge our meta-learners?
- Paired open-ended trailblazer (POET): evolves a parameterized environment  $\Theta_E$  for agent  $\Theta_A$ 
  - Select agents that can solve challenges AND evolve environments so they are solvable

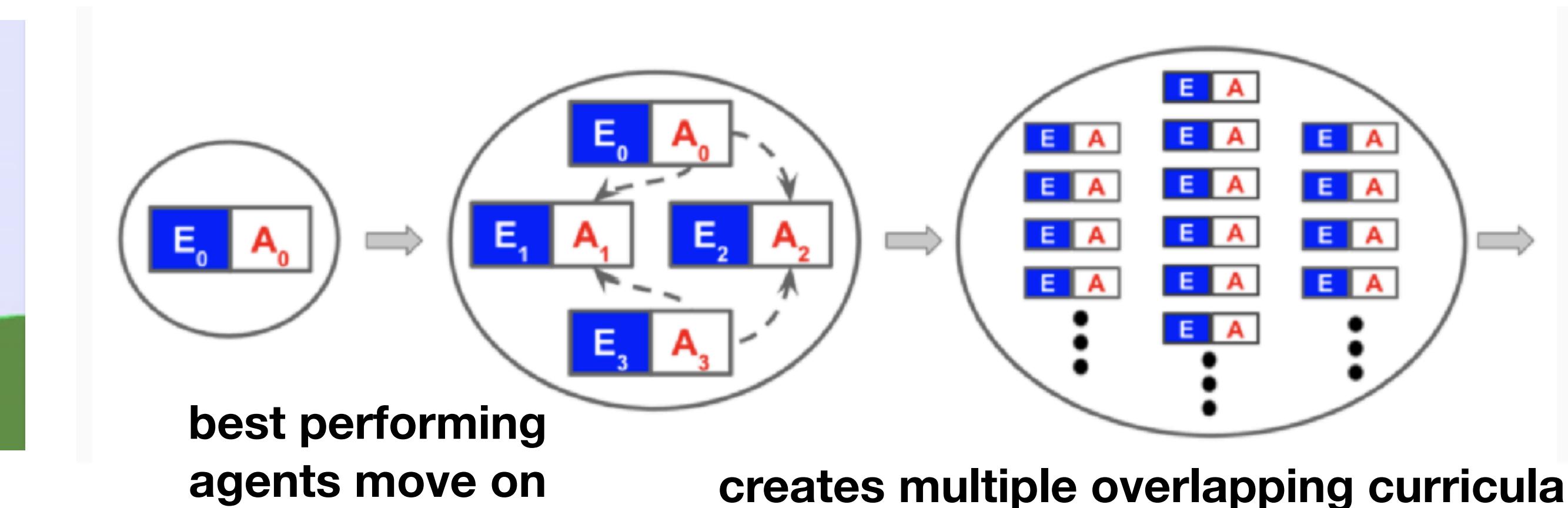
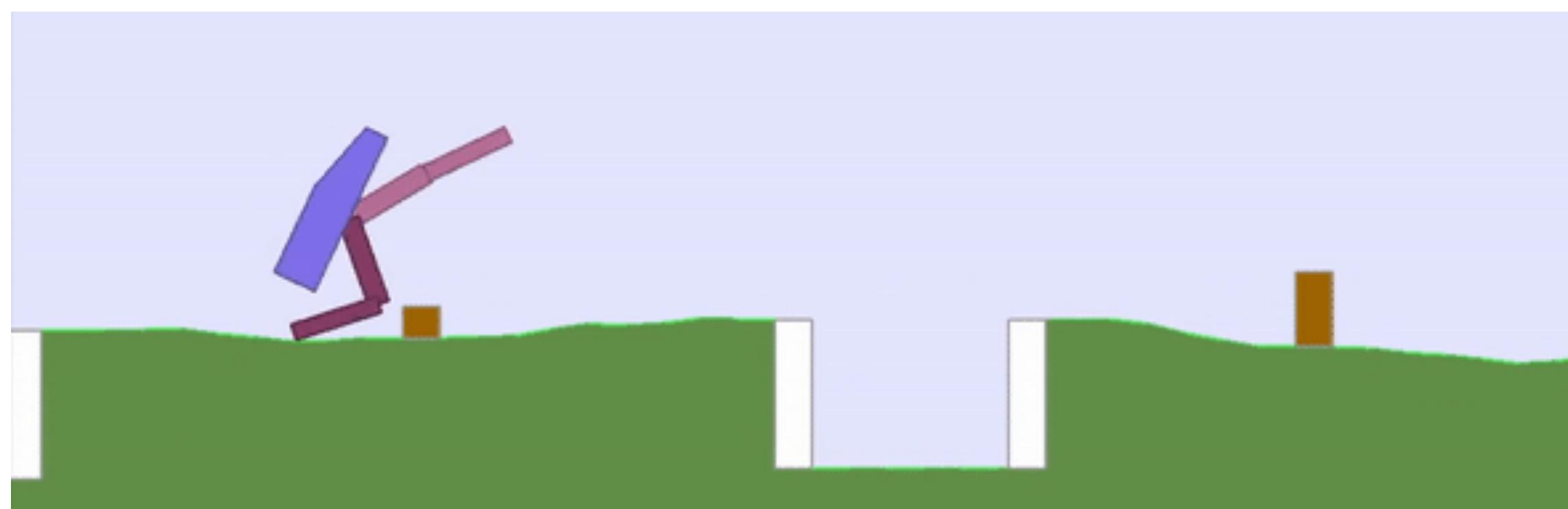


Figure source: Wang et al. 2019

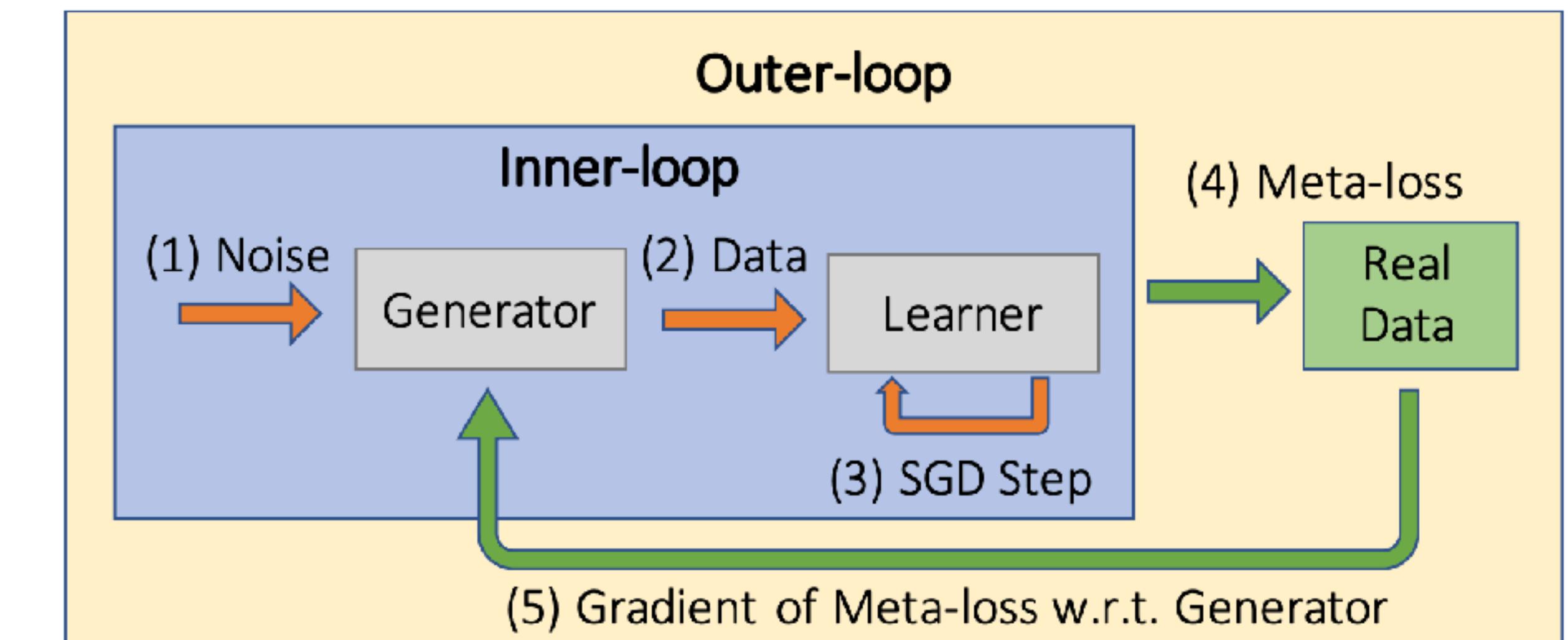
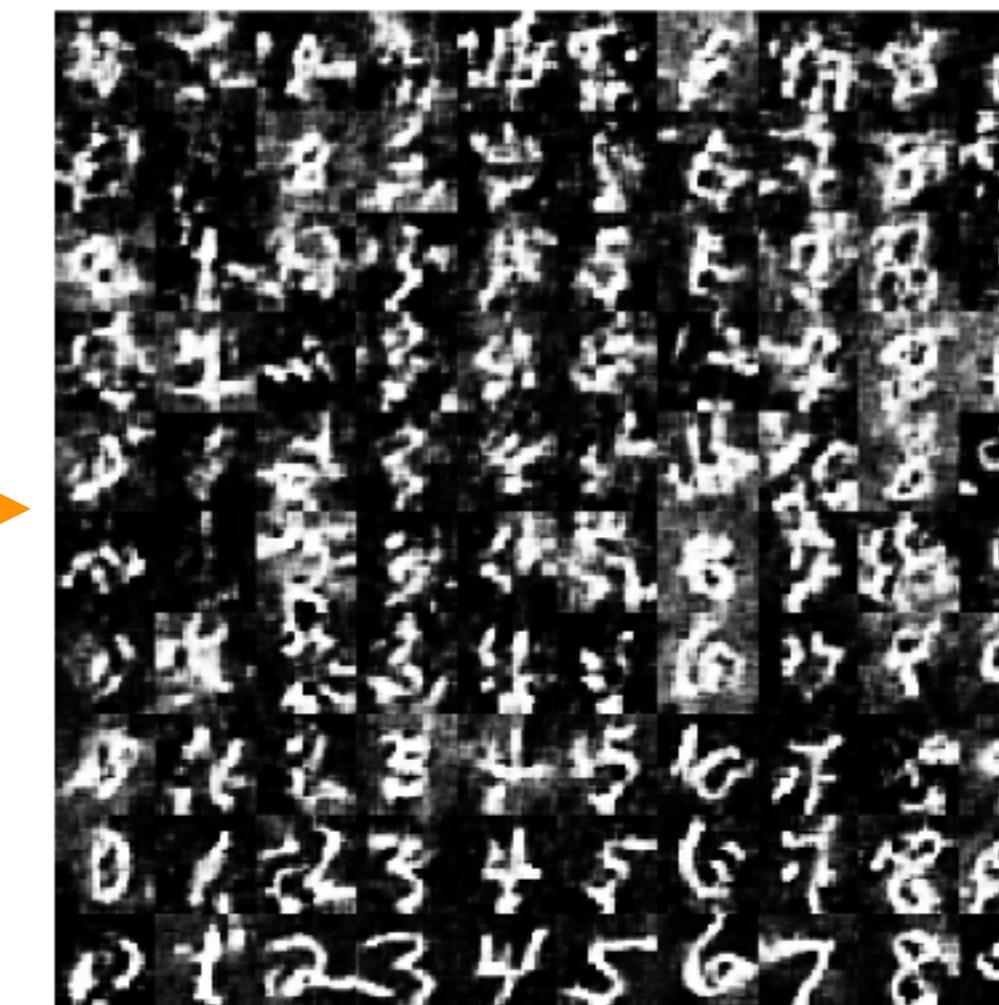
# Does POET scale? Increasingly difficult 3D terrain, 18 degrees of freedom.



# Generative Teaching Networks

- Based on an existing dataset, generate synthetic training data for more efficient training
  - Like dataset distillation, but uses meta-learning to update the generator model
- While POET has limited expressivity (limited to  $\Theta_E$ ), GTNs could produce all sorts of training datasets and environments
  - While being careful not to generate noise: needs some grounding in reality

3 4 2 1 9 5 6 2 1 8  
 8 9 1 2 5 0 0 6 6 4  
 6 7 0 1 6 3 6 3 7 0  
 3 7 7 9 4 6 6 1 8 2  
 2 9 3 4 3 9 8 7 2 5  
 1 5 9 8 3 6 5 7 2 3  
 9 3 1 9 1 5 8 0 8 4  
 5 6 2 6 8 5 8 8 9 9  
 3 7 7 0 9 4 8 5 4 3  
 7 9 6 4 1 0 6 9 2 3

# Thank you

