

Meta-Learning

ACDL 2021, part 1

Joaquin Vanschoren (TU Eindhoven)

Intro: humans can easily learn from a single example



Intro: humans can easily learn from a single example



Intro: humans can easily learn from a single example

Canna Indica 'Picasso'



Intro: humans can easily learn from a single example

Canna Indica 'Picasso'

?



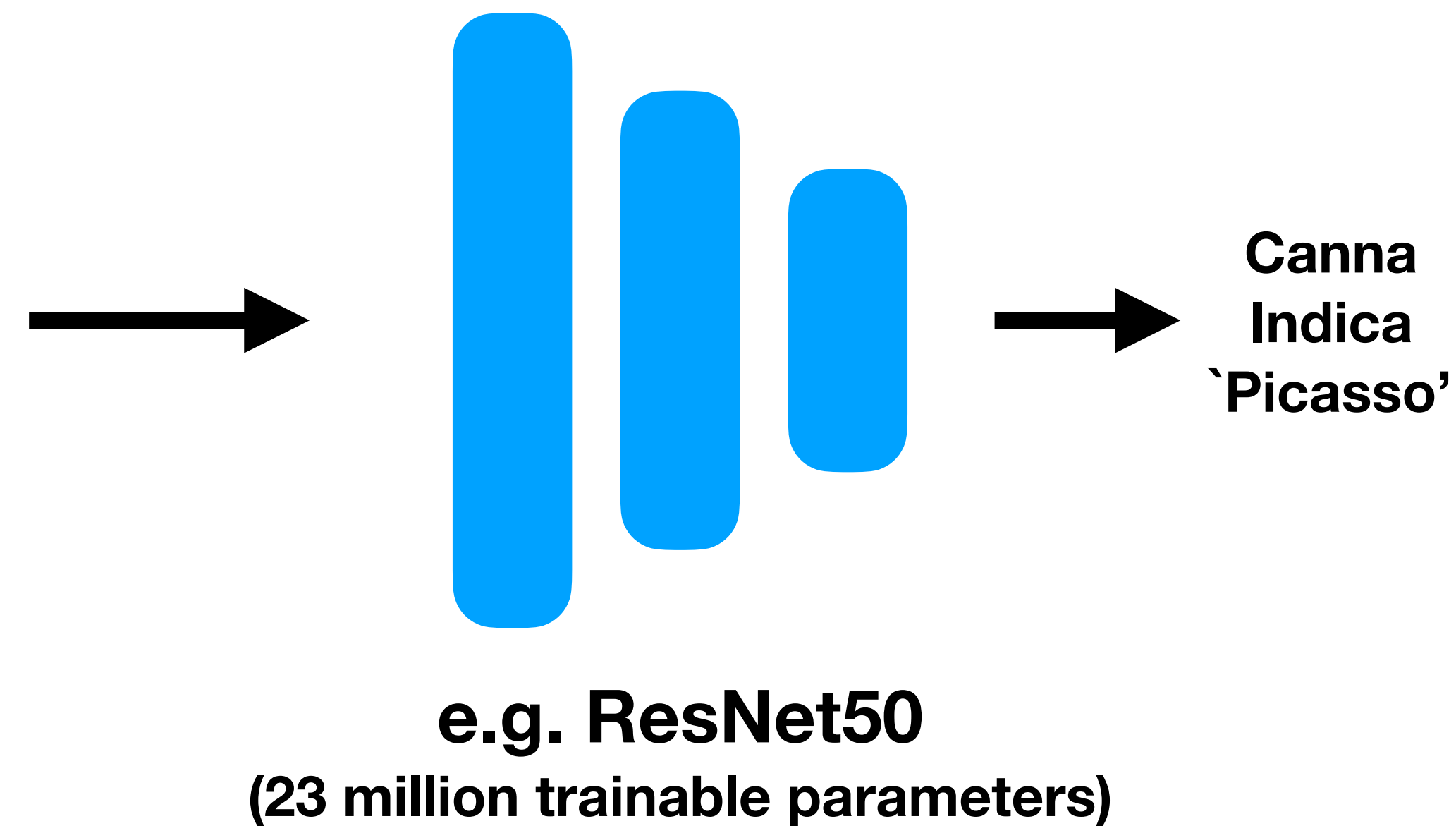
later
→



train

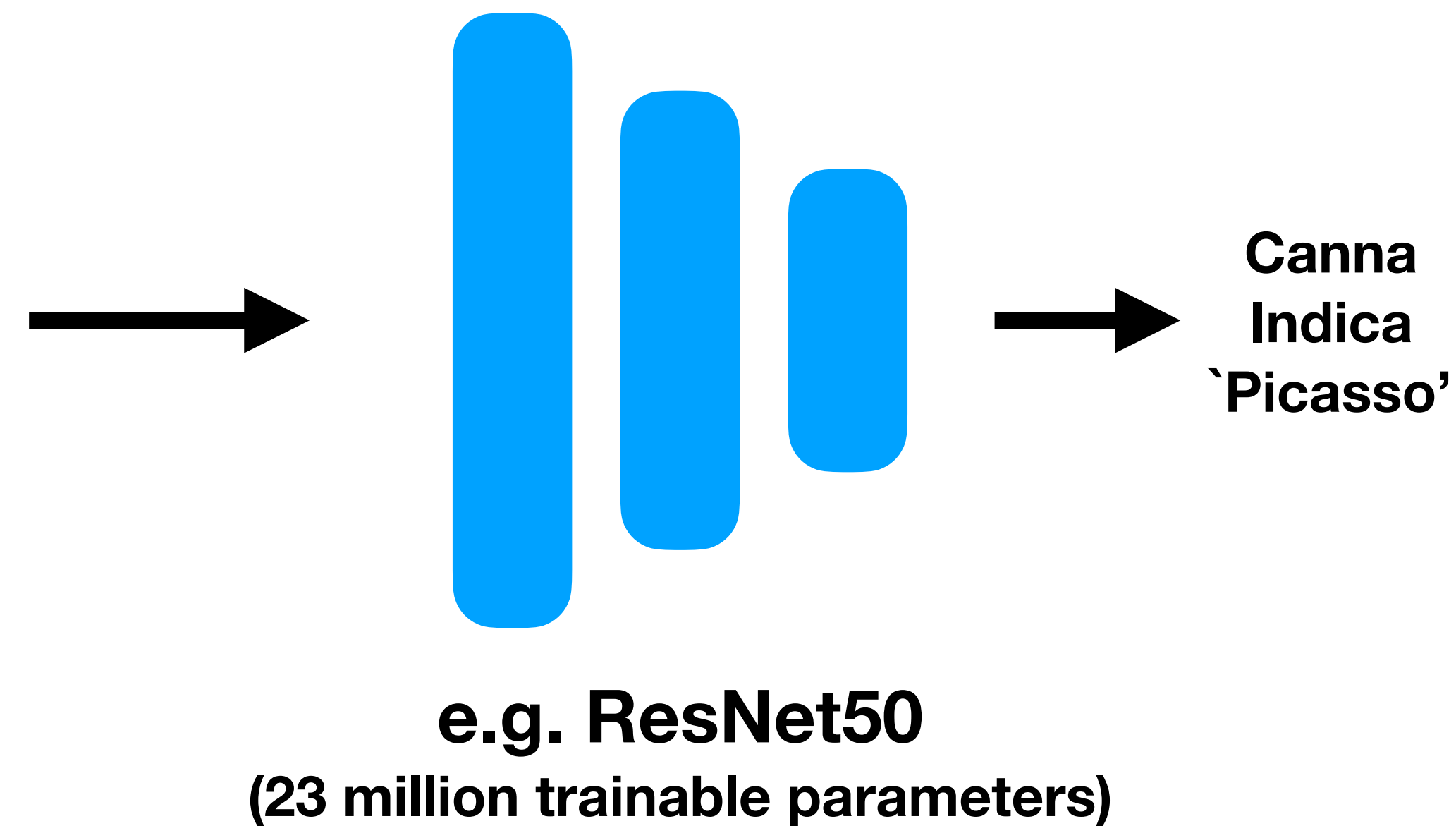
test

Can a computer learn from a single example?



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

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

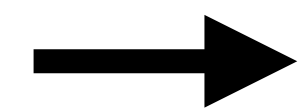


ImageNet
(14 million images)

Pretrain



e.g. ResNet50
(23 million trainable parameters)



Canna
Indica
'Picasso'

Transfer learning?

Target task



Source task

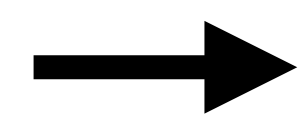


ImageNet
(14 million images)

↓ Pretrain



e.g. ResNet50
(23 million trainable parameters)



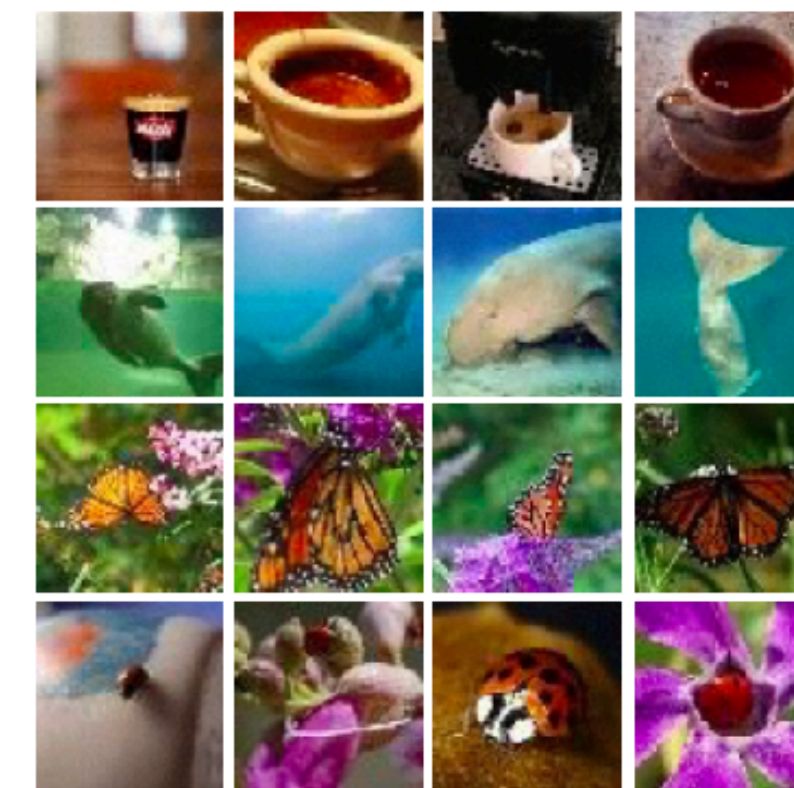
Canna
Indica
'Picasso'

Transfer learning?

Target task



Source task



ImageNet
(14 million images)

Pretrain

Finetune



e.g. ResNet50
(23 million trainable parameters)

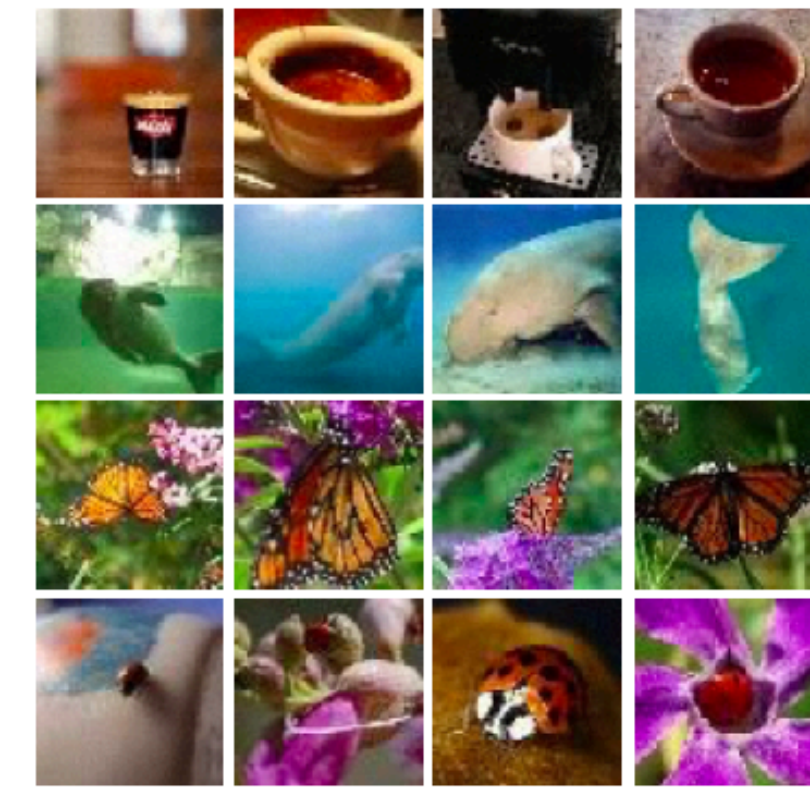
Canna
Indica
'Picasso'

Transfer learning?

Target task



Source task



ImageNet
(14 million images)

Pretrain

Finetune



Canna
Indica
'Picasso'

e.g. ResNet50
(23 million trainable parameters)

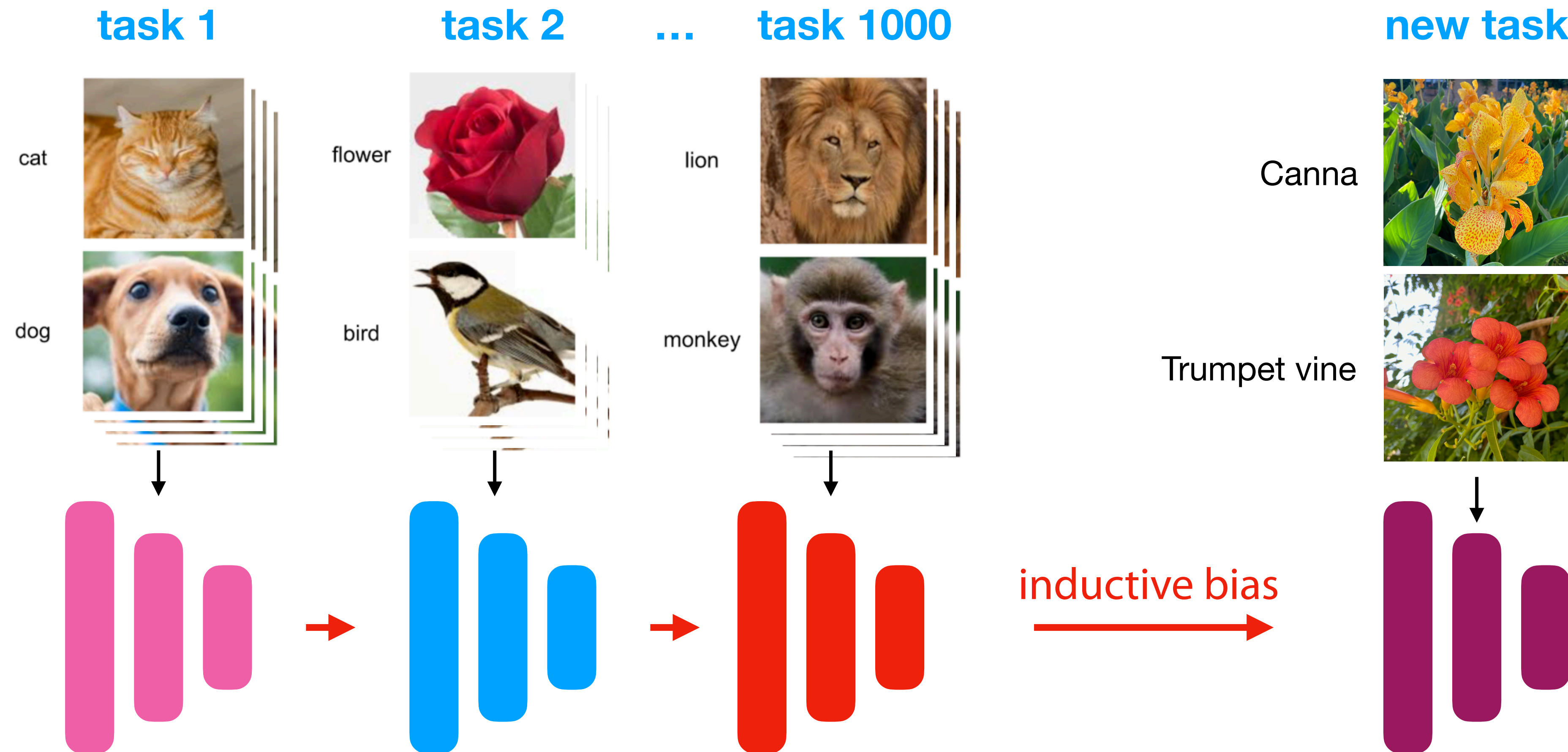
A single source task (e.g. ImageNet) may not generalize well to the test task

Meta-learning

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

Learn what works well, and learn how to transfer that to new tasks

Prepare yourself to learn new things faster

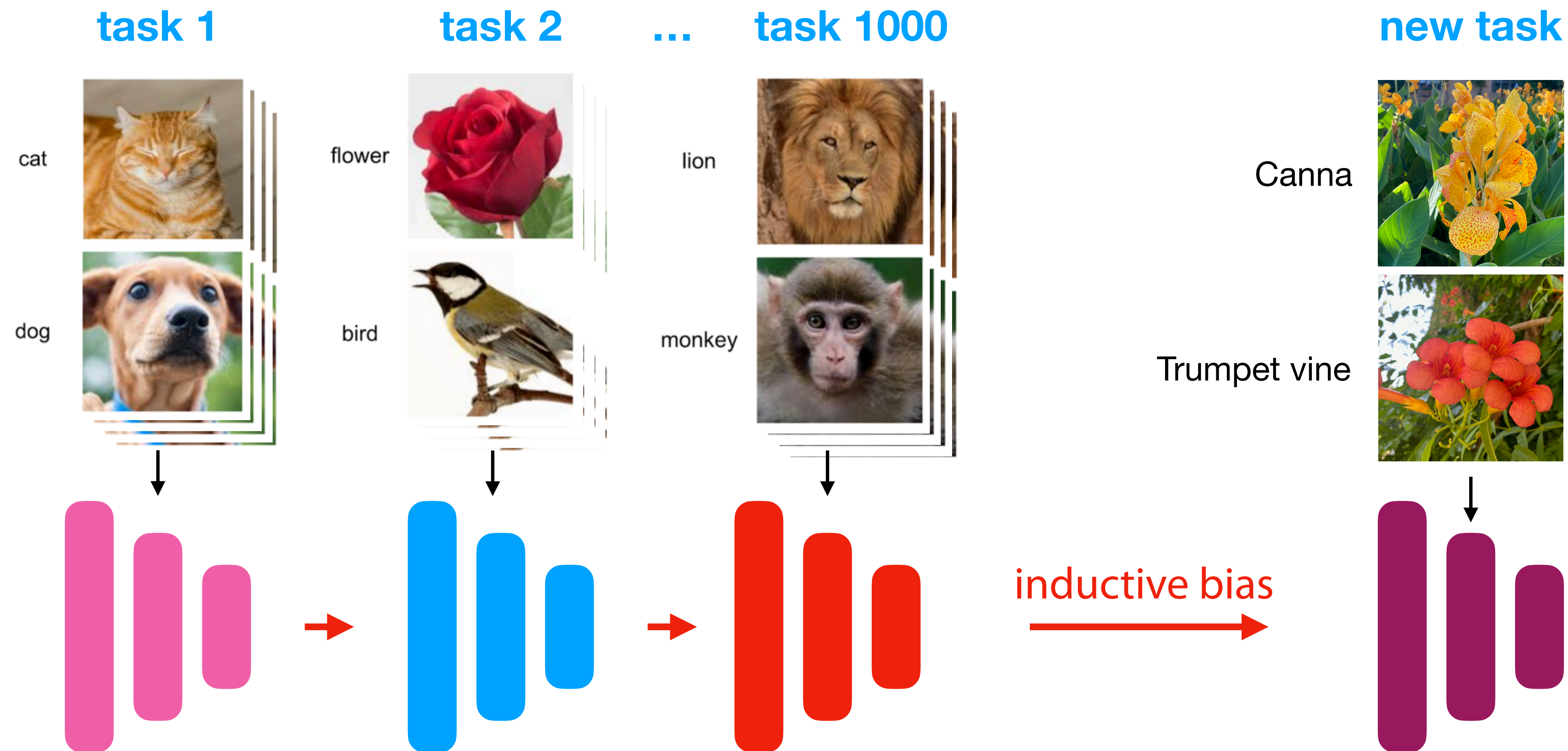


Meta-learning

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

Learn what works well, and learn how to transfer that to new tasks

Prepare yourself to learn new things faster



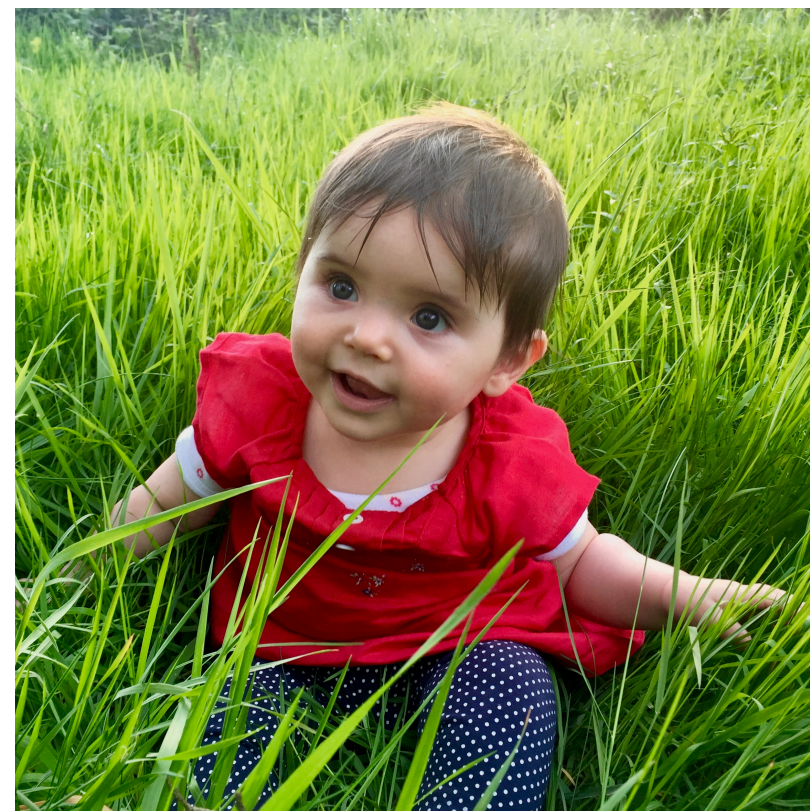
Useful in many real-life situations: rare events, test-time constraints, data collection costs, privacy issues,...

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

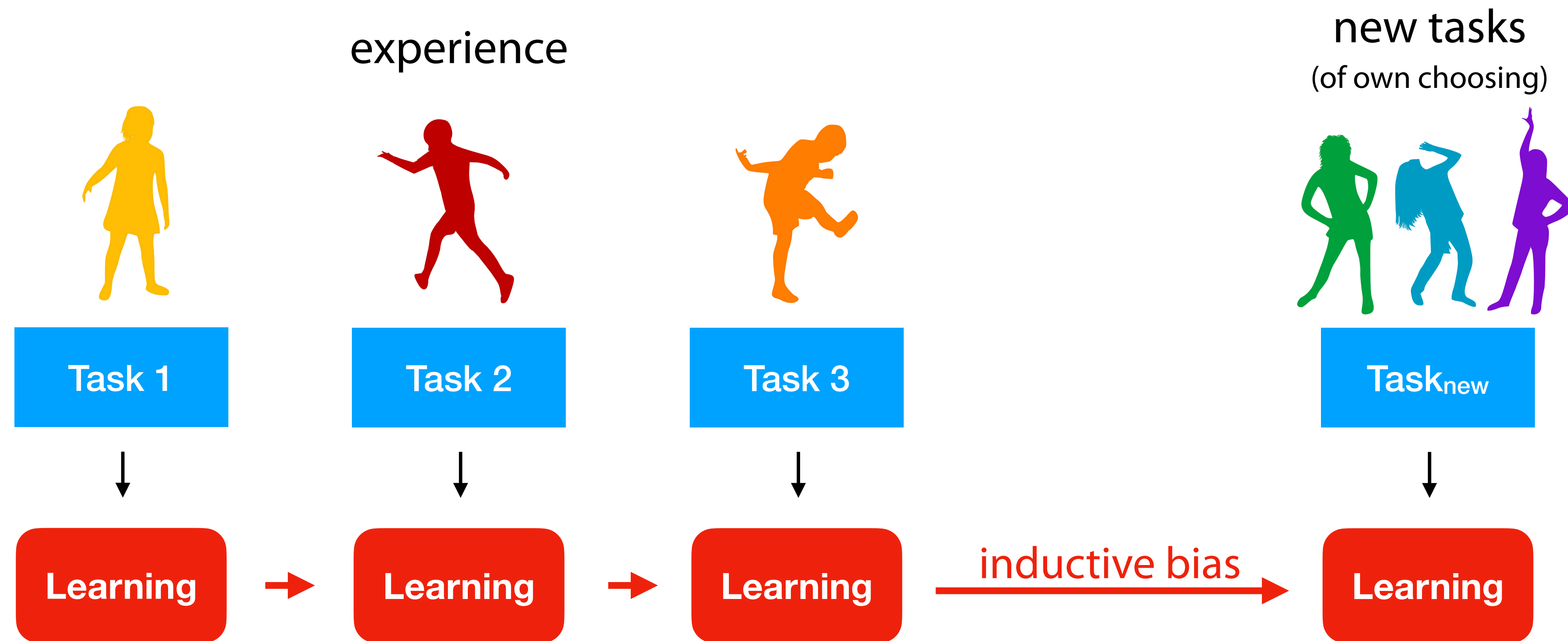


year 4



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

learning compositionality (e.g. language)

easy enough for humans, hard for computers

Training

-  dax
-  zup
-  lug
-  wif

-    lug blicket wif
-    wif blicket dax
-   lug kiki wif
-   wif kiki dax

Test





dax blicket zup?











wif blicket dax kiki lug?

learning compositionality (e.g. language)

easy enough for humans, hard for computers

Training

-  dax
-  zup
-  lug
-  wif

-    lug blicket wif
-    wif blicket dax
-   lug kiki wif
-   wif kiki dax

Test

   dax blicket zup?

wif blicket dax kiki lug?

learning compositionality (e.g. language)

easy enough for humans, hard for computers

Training

-  dax
-  zup
-  lug
-  wif

   lug blicket wif

   wif blicket dax

  lug kiki wif

  wif kiki dax

Test

   dax blicket zup?

    wif blicket dax kiki lug?

learning compositionality (e.g. language)

easy enough for humans, hard for computers

Training

-  dax
-  zup
-  lug
-  wif

   lug blicket wif

   wif blicket dax

  lug kiki wif

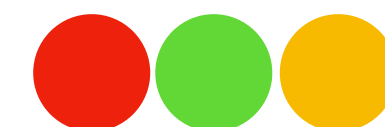
  wif kiki dax

Test

   dax blicket zup?

    wif blicket dax kiki lug?

Common mistakes



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

learning compositionality (e.g. language)

easy enough for humans, hard for computers

Training

-  dax
-  zup
-  lug
-  wif

   lug blicket wif

   wif blicket dax

  lug kiki wif


  wif kiki dax





Test

   dax blicket zup?

    wif blicket dax kiki lug?

Common mistakes

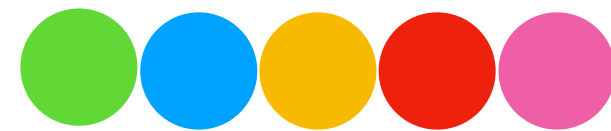
   one-to-one bias:
assume that every word is one color

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

What if there is *no* training data?

humans can still solve problems by making assumptions

Item pool



Test

zup?

zup zup?

dax zup?

zup tufa?

zup wif zup?

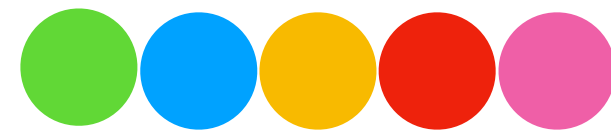
zup wif blicket?

blicket wif zup?

What if there is *no* training data?

humans can still solve problems by making assumptions

Item pool



Test



zup?

zup zup?

dax zup?

zup tufa?

zup wif zup?

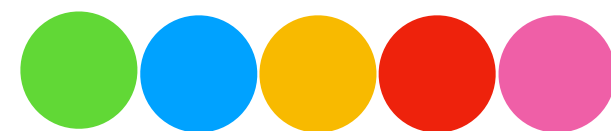
zup wif blicket?

blicket wif zup?

What if there is *no* training data?

humans can still solve problems by making assumptions

Item pool



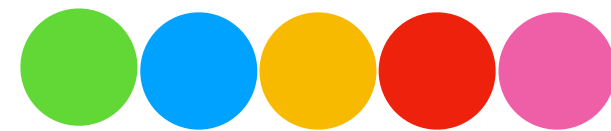
Test

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






What if there is *no* training data?

humans can still solve problems by making assumptions

Item pool



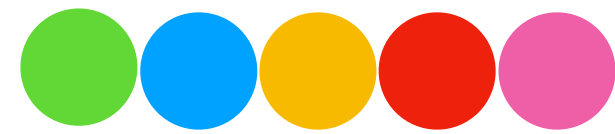
Test

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

What if there is *no* training data?

humans can still solve problems by making assumptions

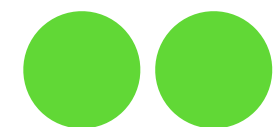
Item pool



Test



zup?



zup zup?



dax zup?



zup tufa?

zup wif zup?

zup wif blicket?

blicket wif zup?



What if there is *no* training data?

humans can still solve problems by making assumptions

Item pool



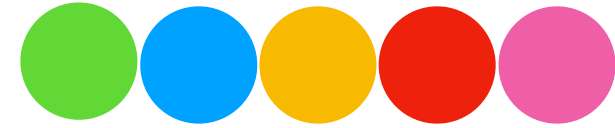
Test

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








What if there is *no* training data?

humans can still solve problems by making assumptions

Item pool



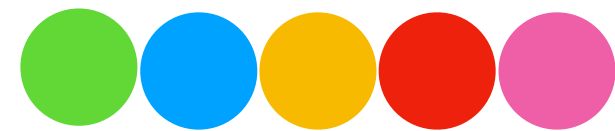
Test

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



What if there is *no* training data?

humans can still solve problems by making assumptions

Item pool



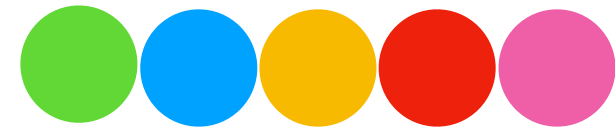
Test

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






What if there is *no* training data?

humans can still solve problems by making assumptions

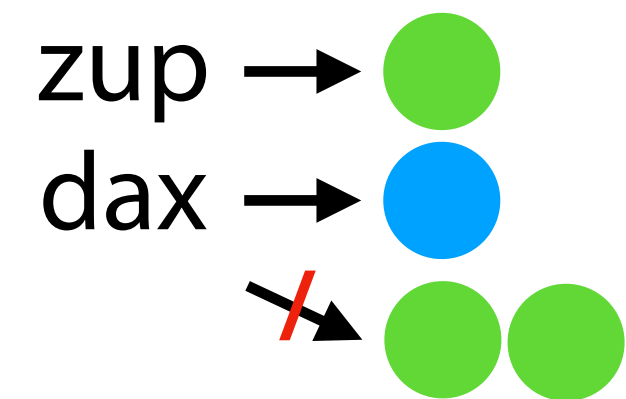
Item pool



Test

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

Commonly used assumptions:



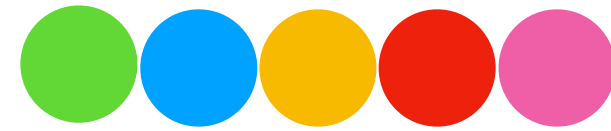
one-to-one bias:

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



What if there is *no* training data?

humans 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:


zup → 

dax → 

~~→~~ 

one-to-one bias:

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

dax zup → 

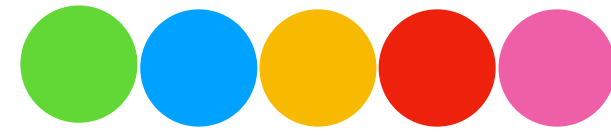
concatenation bias:

assume that order is always left-to-right







What if there is *no* training data?

humans can still solve problems by making assumptions

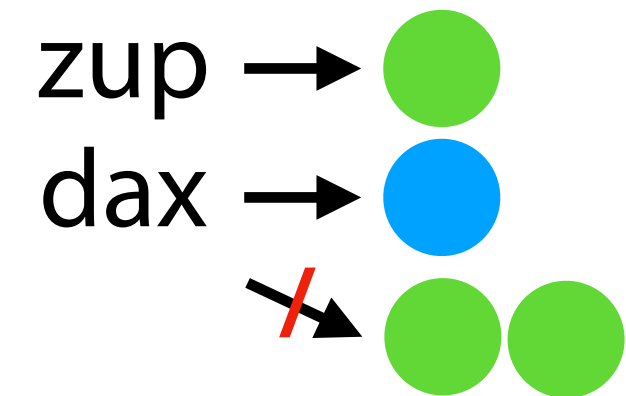
Item pool



Test

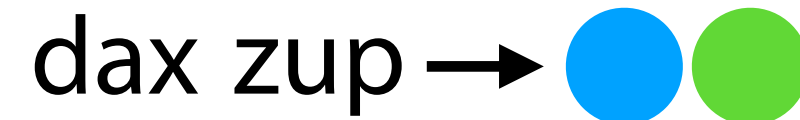
-  zup?
-  zup zup?
-  dax zup?
-  zup tufa?
-  zup wif zup?
-  zup wif blicket?
-  blicket wif zup?

Commonly used assumptions:



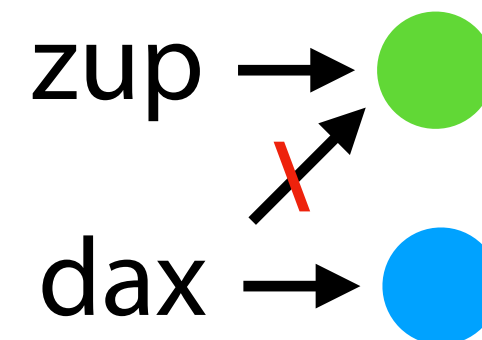
one-to-one bias:

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



concatenation bias:

assume that order is always left-to-right



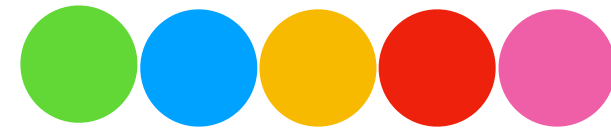
mutual exclusivity:

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






What if there is *no* training data?

humans can still solve problems by making assumptions

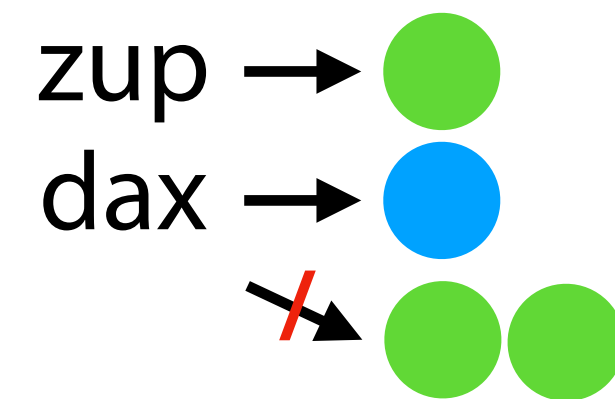
Item pool



Test

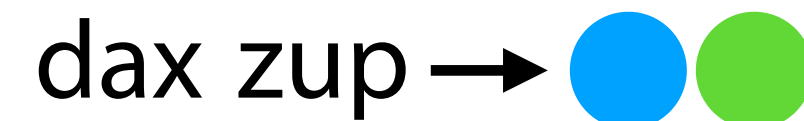
-  zup?
-  zup zup?
-  dax zup?
-  zup tufa?
-  zup wif zup?
-  zup wif blicket?
-  blicket wif zup?

Commonly used assumptions:



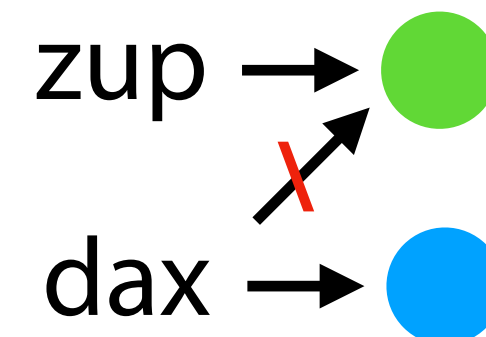
one-to-one bias:

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:

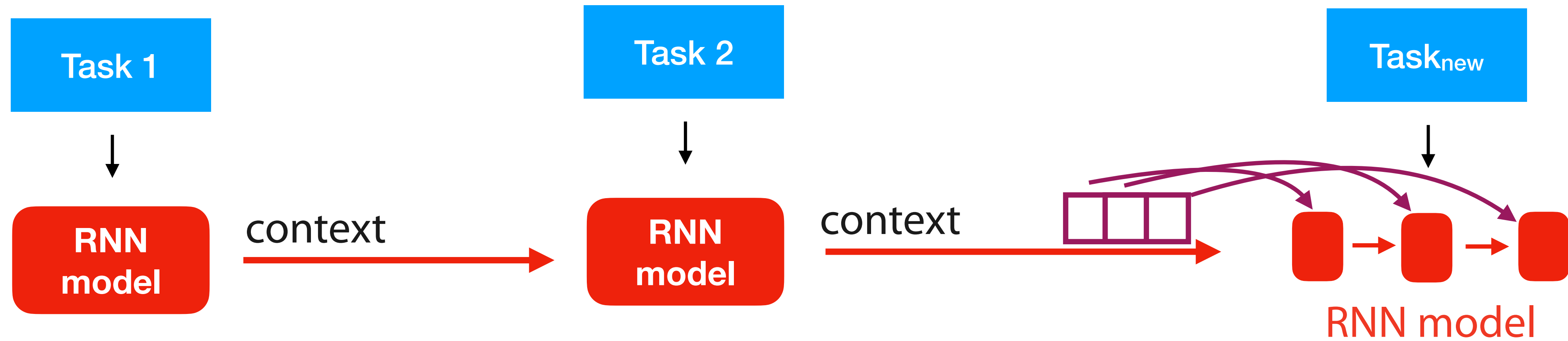
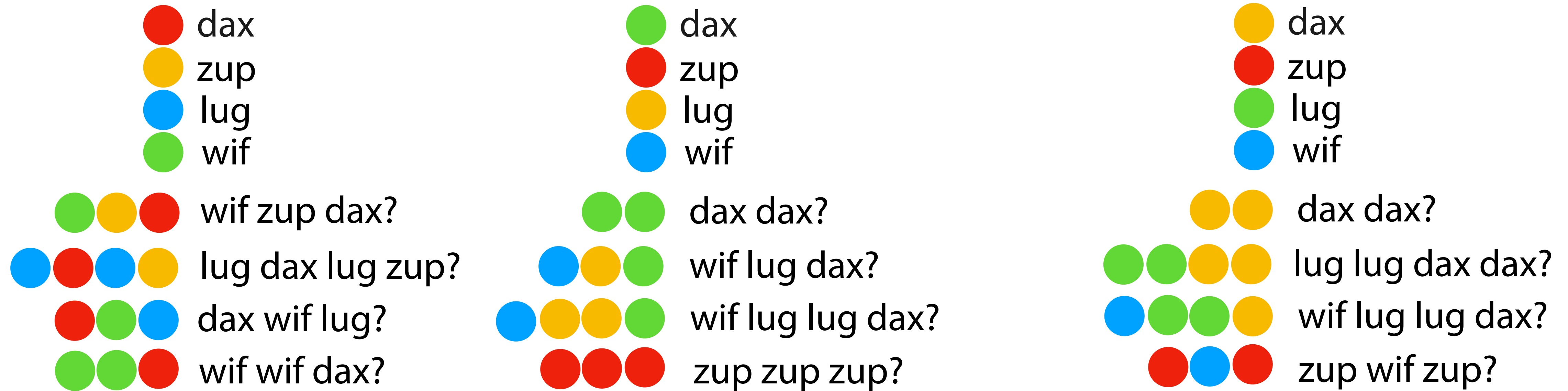
if object has a name, it doesn't need another name

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

These assumptions (inductive biases) are necessary for learning

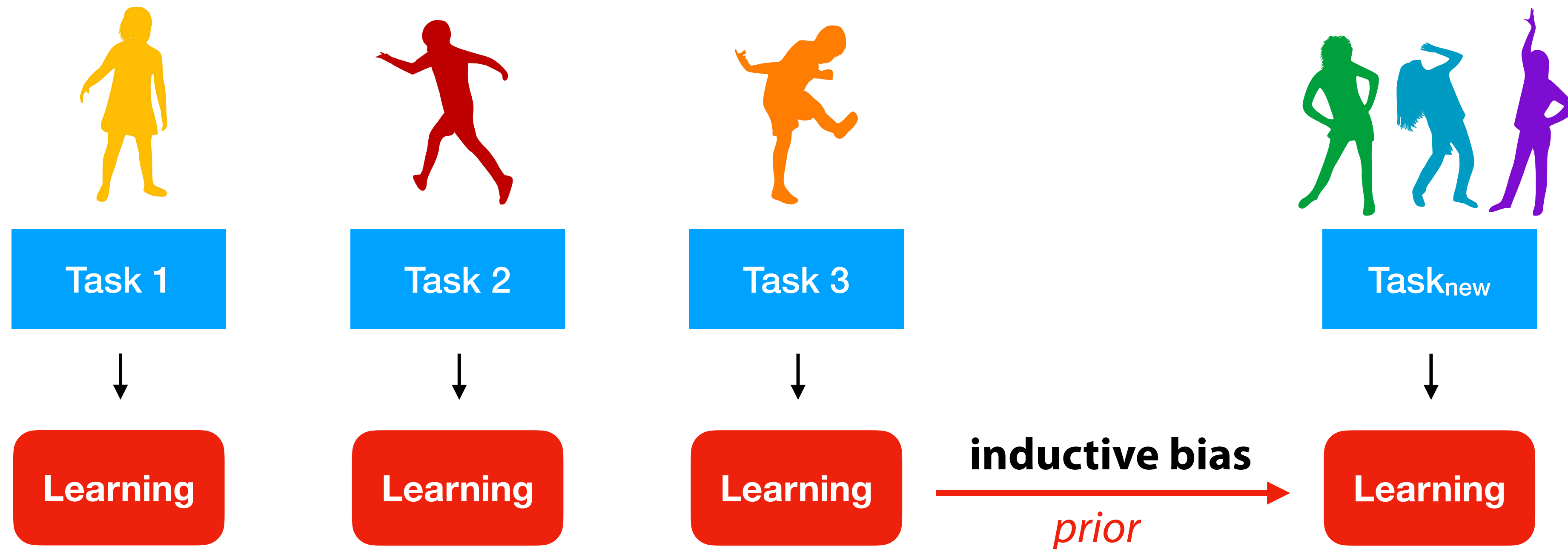
Meta-learning inductive biases

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 **starting from random weights**, **learn good initial weights**
- Instead of **normal backpropagation**, **learn how to better update the weights**
- Instead of **standard loss/reward function**, **learn a better loss/reward function**
- Instead of **off-the-shelf model architectures**, **learn better architectures (and hyperparameter)**

Terminology

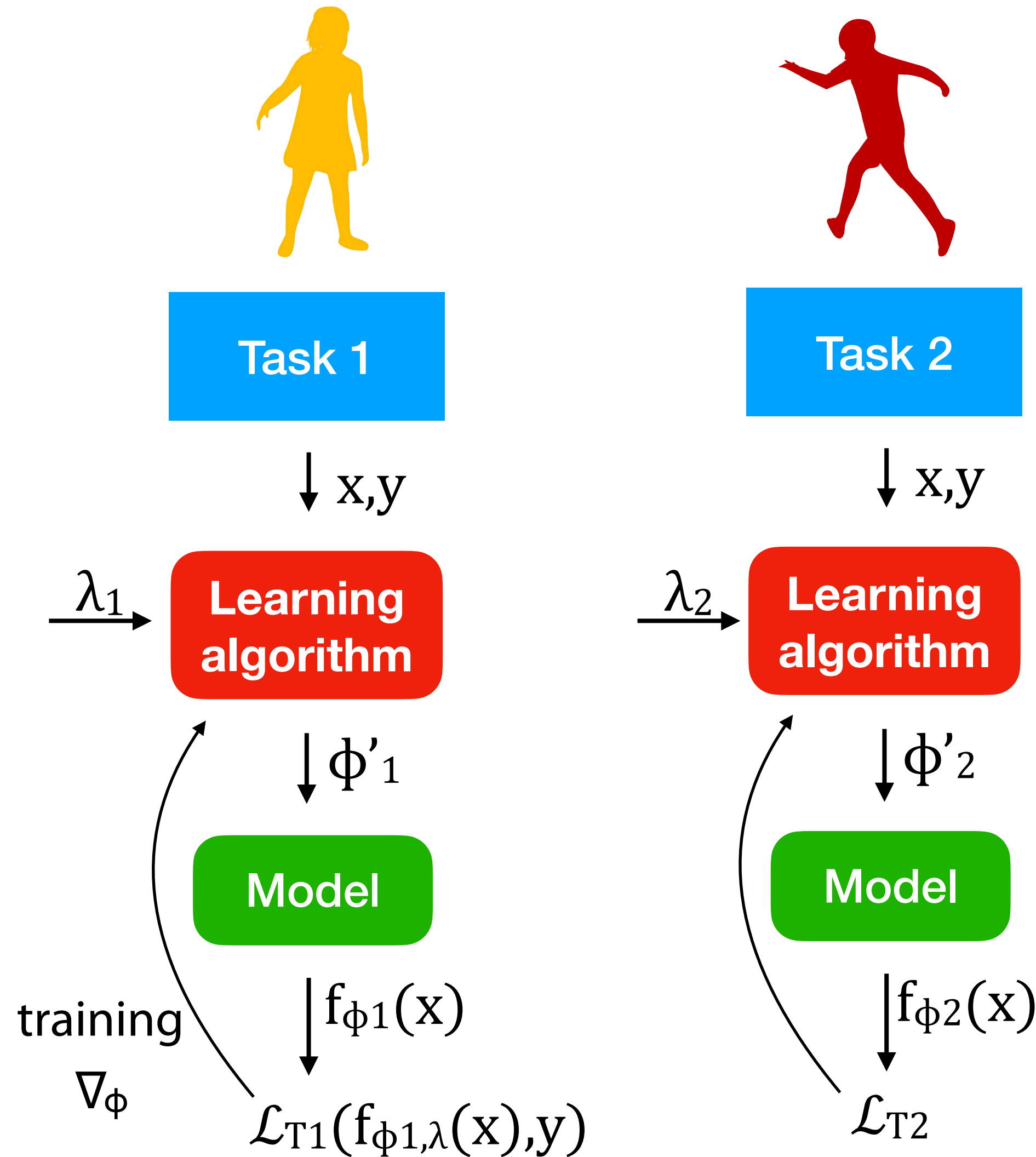
Task: distribution of samples $q(x)$
outputs y , loss $\mathcal{L}(x,y)$

Learner: model parameters ϕ ,
hyper-parameters λ
optimizer

$$\phi^* = \operatorname{argmax}_{\phi} \log p(\phi | T)$$

$$= \operatorname{argmin}_{\phi} \mathcal{L}(f_{\phi,\lambda}(x), y)$$

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

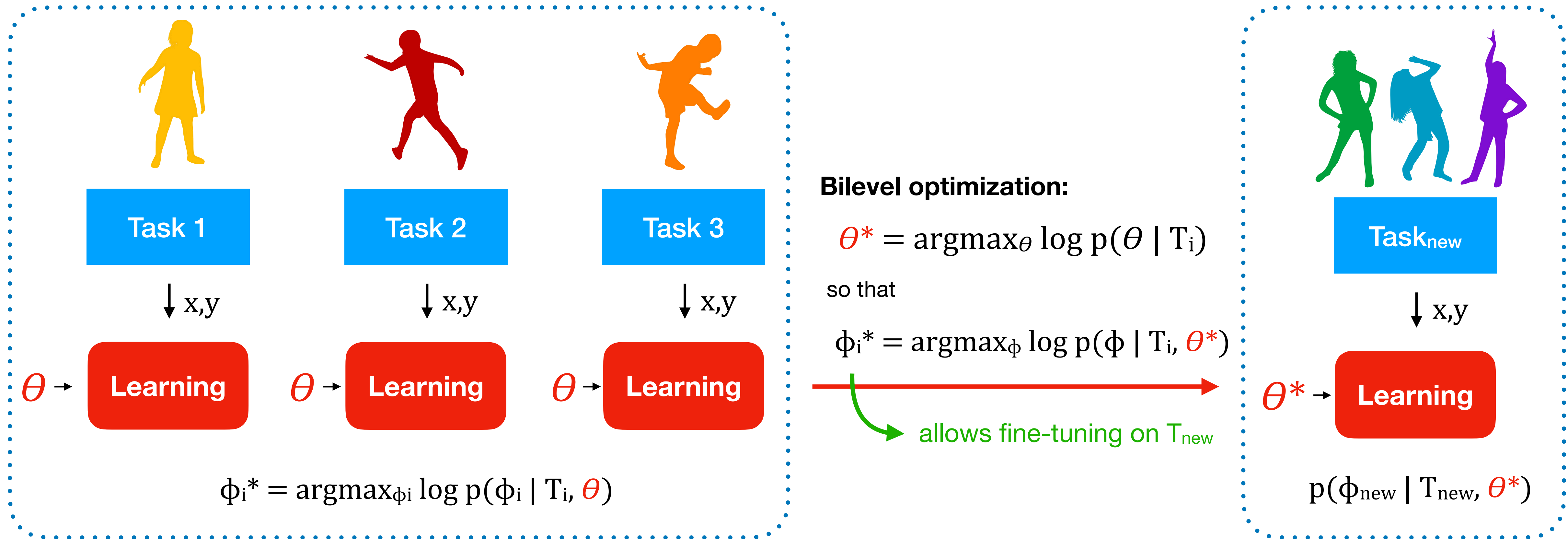


for reinforcement learning:
 $q(x_1) + \text{transition } q(x_{t+1} | x_t, y_t)$
 $\mathcal{L} = (\text{neg}) \text{ reward}$

Strategy 1: bilevel optimization

parameterize some aspect of the learner that we want to learn as meta-parameters θ

meta-learn θ across tasks



θ (prior), could encode an initialization ϕ , the hyperparameters λ , the optimizer,...

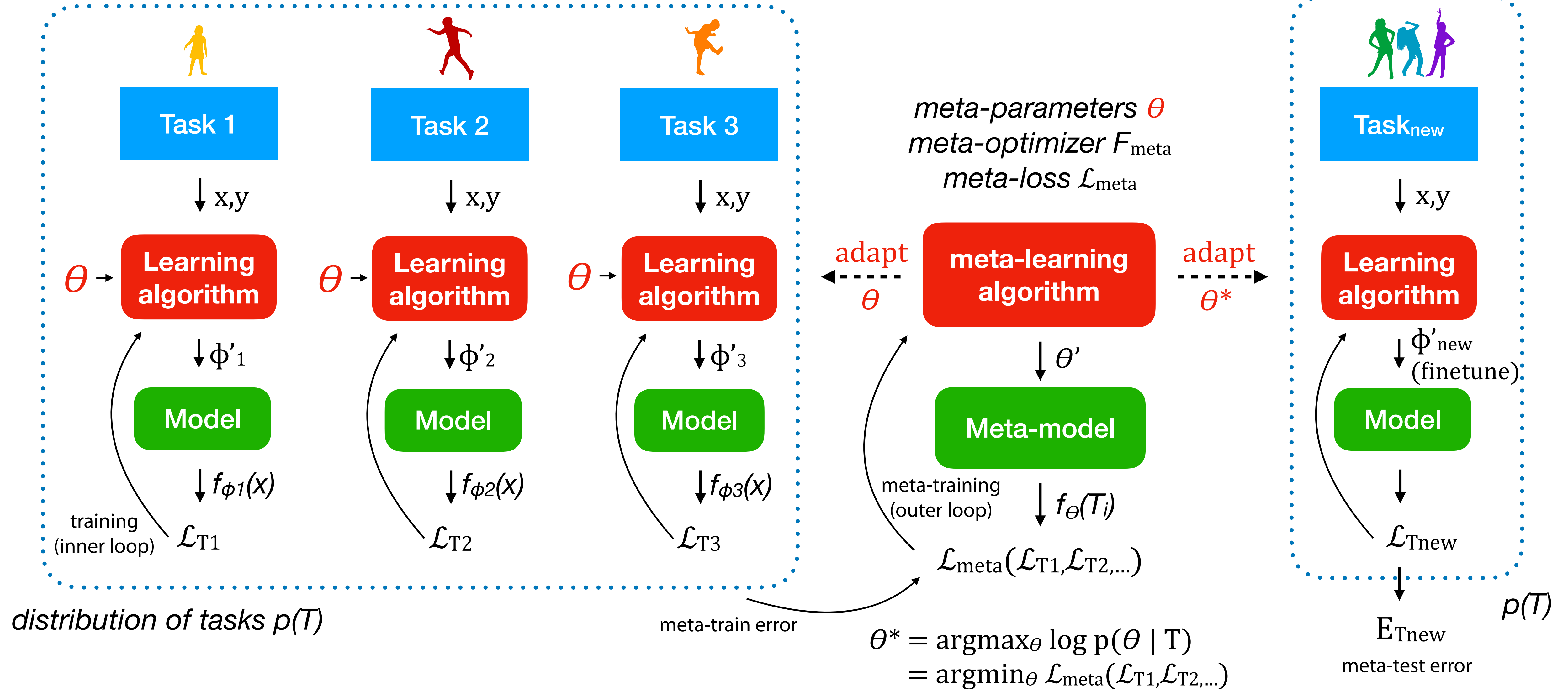
Learned θ^* should *learn* T_{new} from small amount of data, yet *generalize* to a large number of tasks

Meta-learning with bilevel optimization

meta-train

every task is a meta-training example

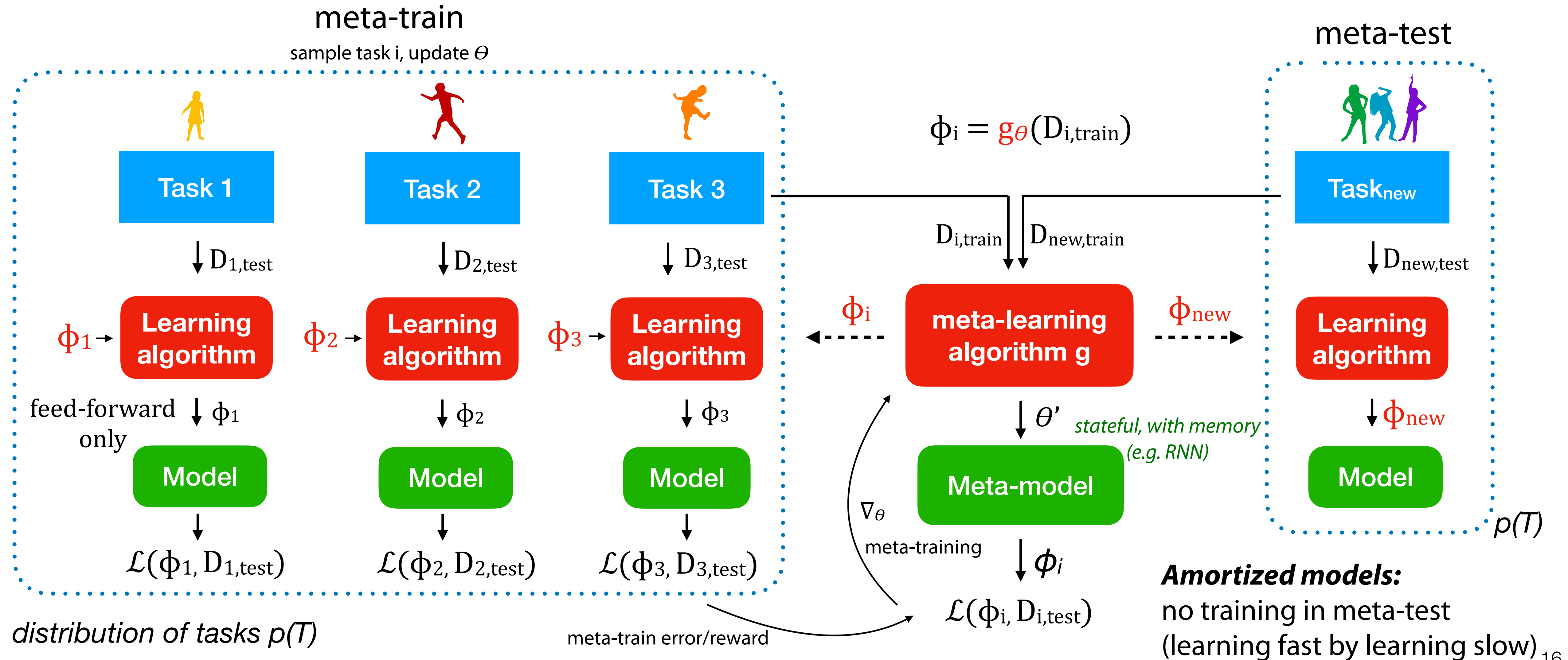
meta-test



Strategy 2: black-box models

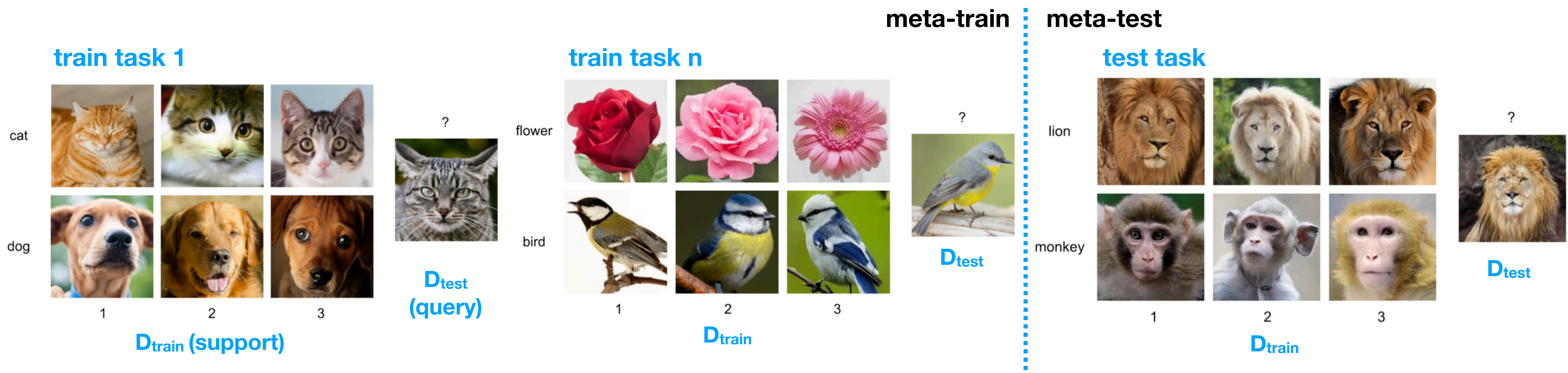
black box meta-model g_θ : predicts ϕ given D_{train} (θ is hidden)

hypernetwork where input embedding learned across tasks

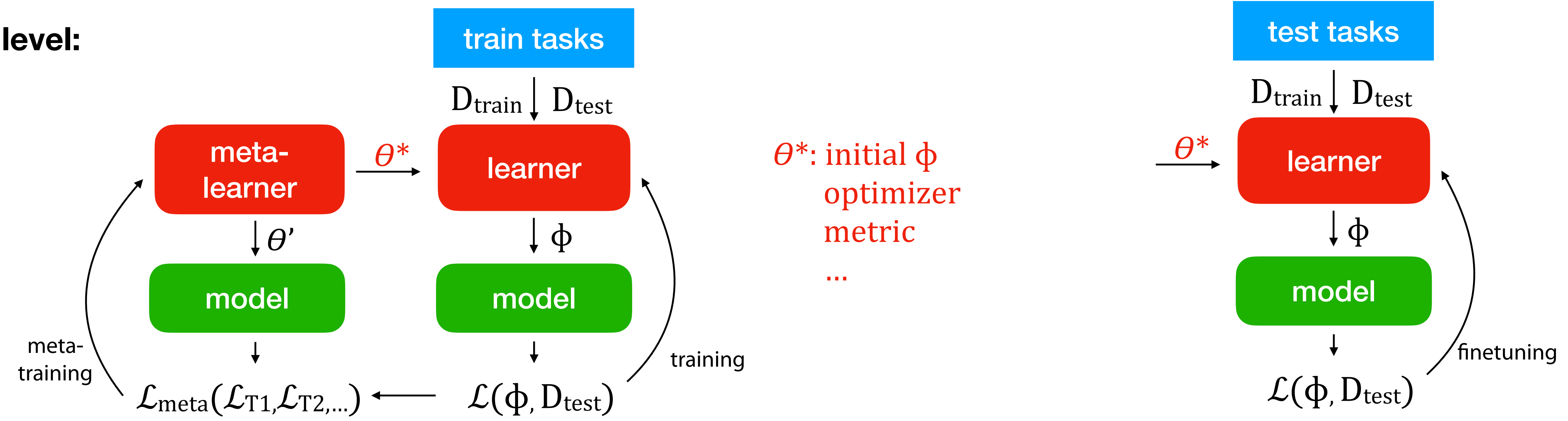


Example: few-shot classification

2-way, 3-shot



Bilevel:



Example: few-shot classification

train task 1



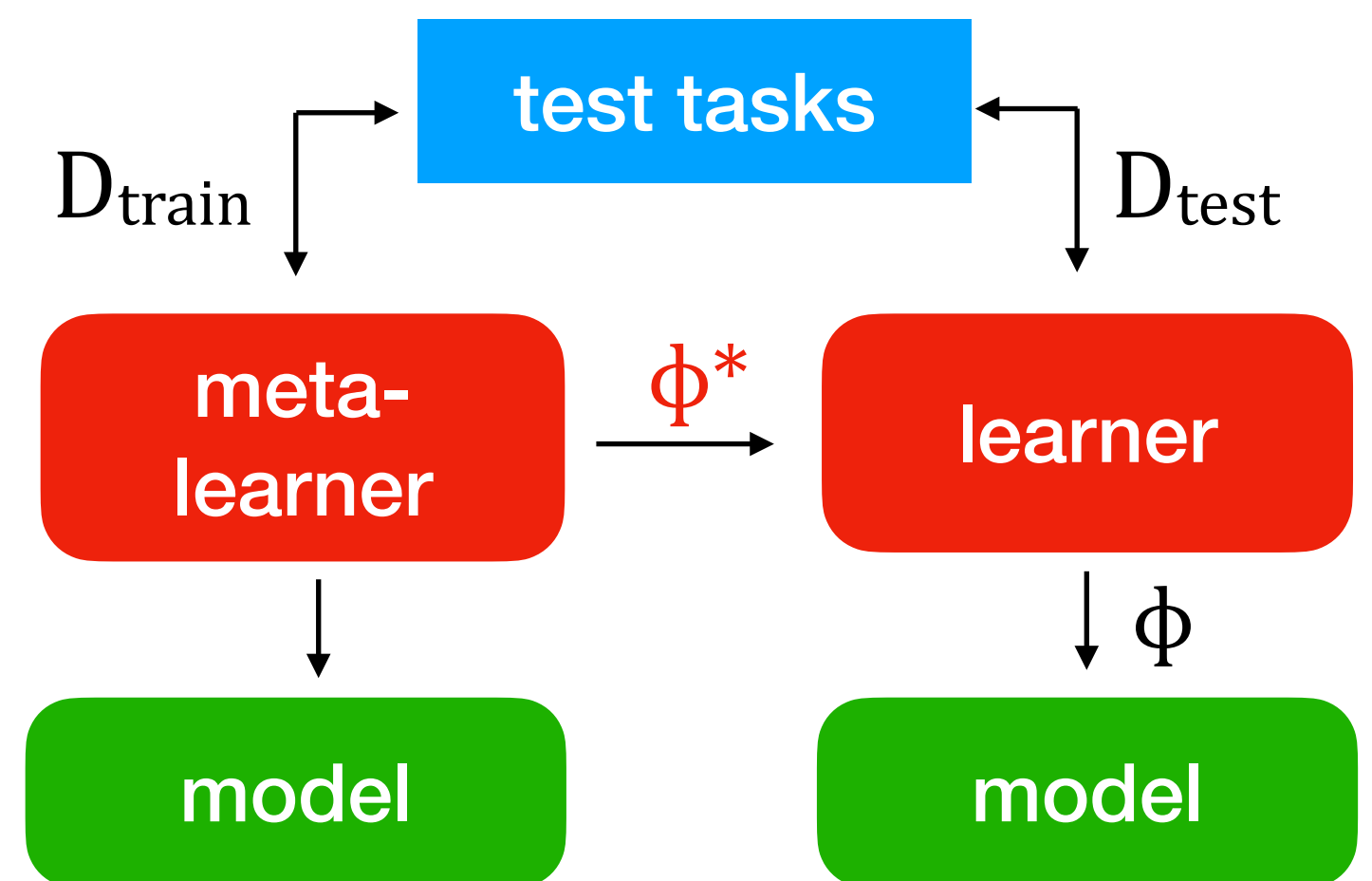
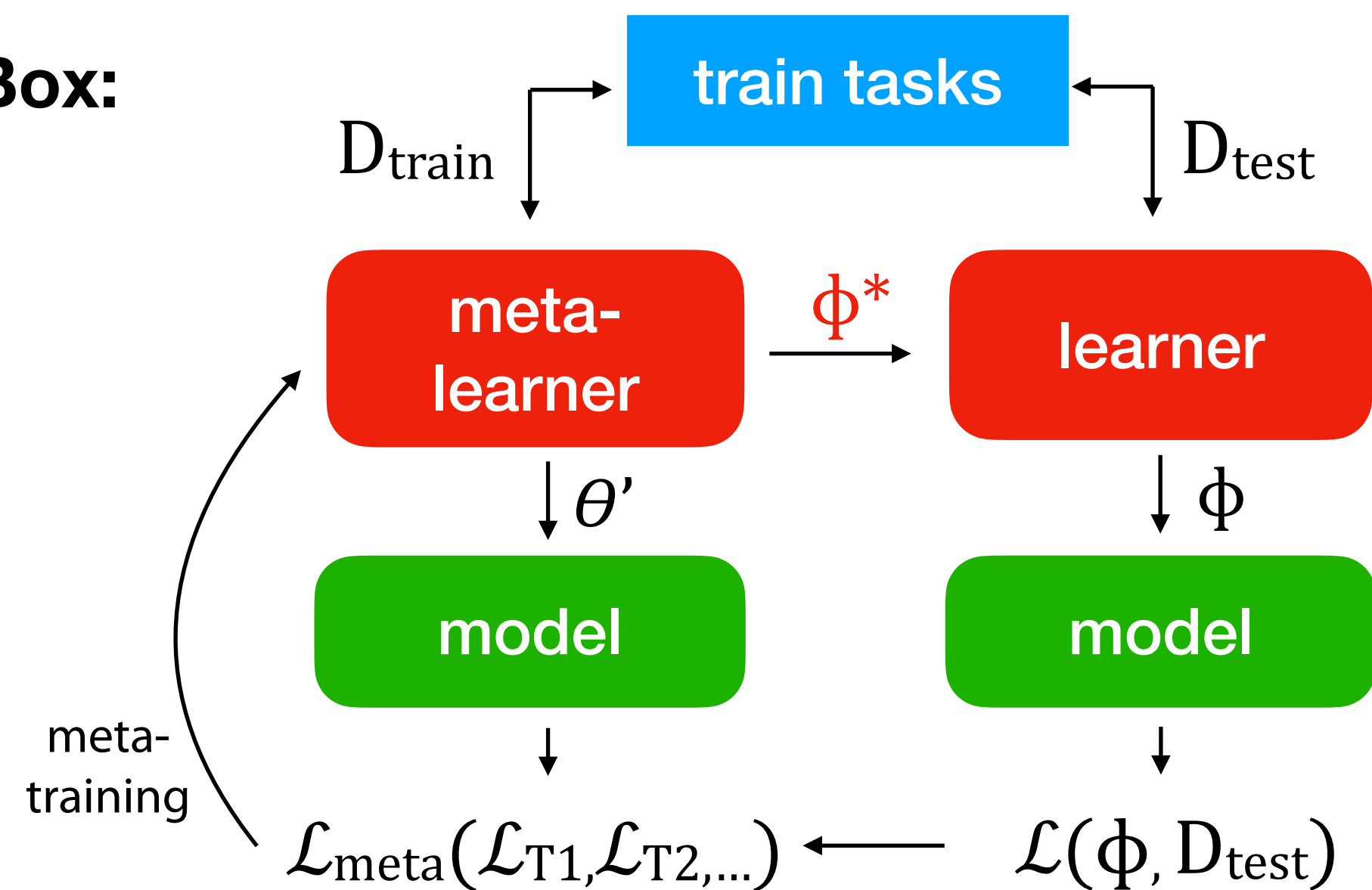
train task n



test task



Black Box:



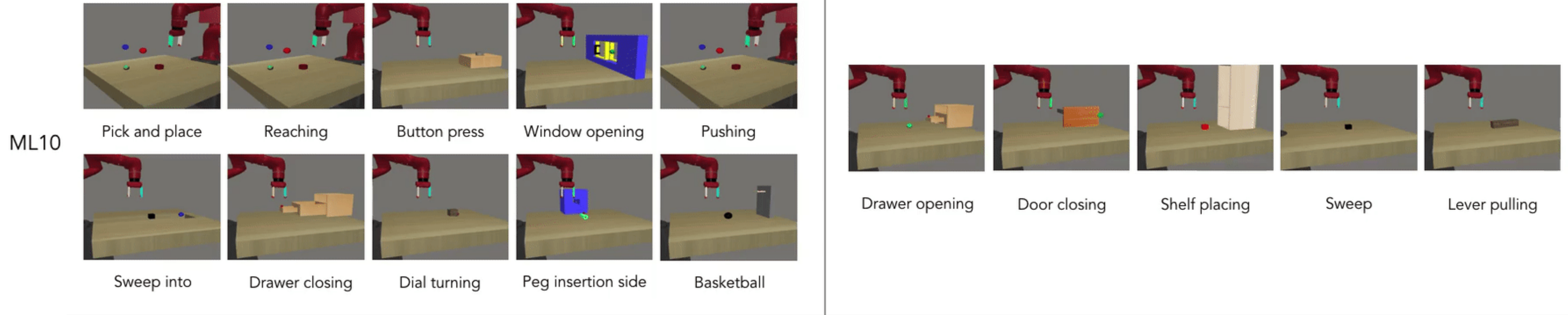
Example: meta-reinforcement learning

initial state:
randomized object
and goals positions

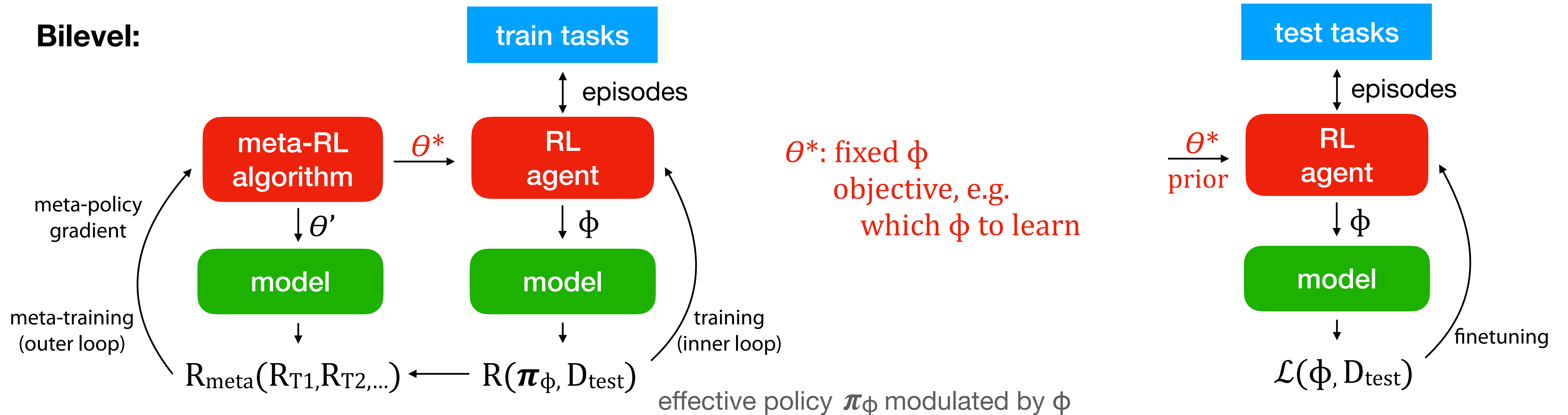
Train tasks

Test tasks

other initial states
or related tasks



Bilevel:



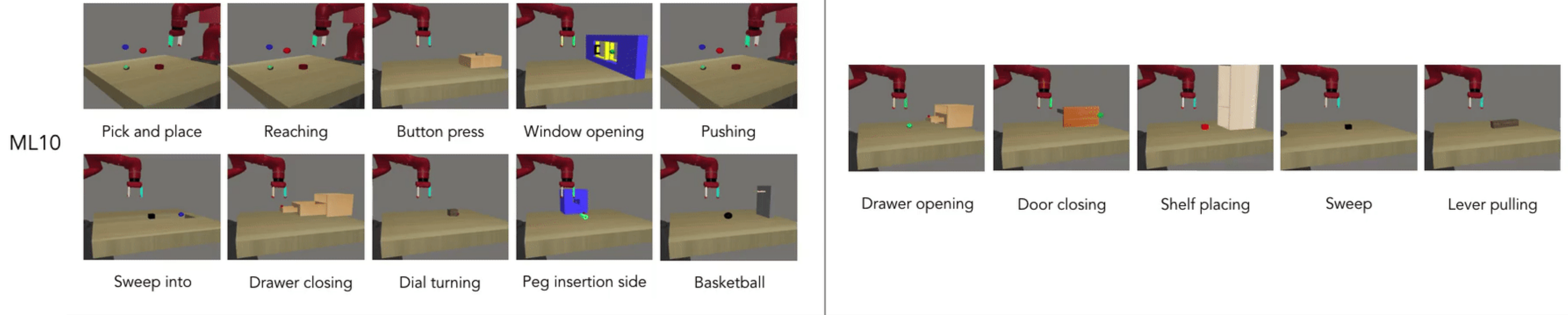
Example: meta-reinforcement learning

initial state:
randomized object
and goals positions

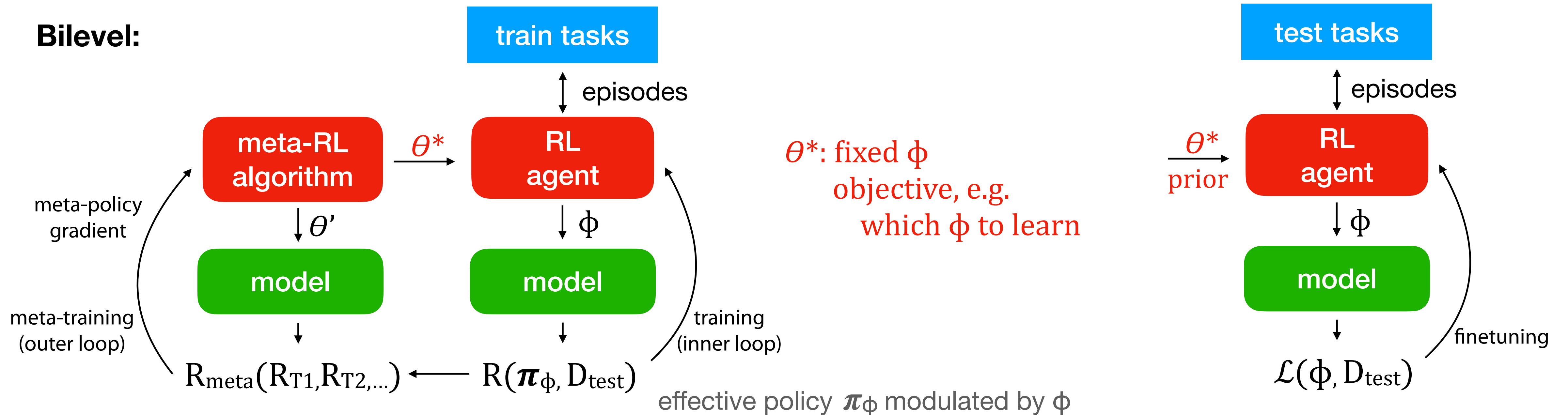
Train tasks

Test tasks

other initial states
or related tasks



Bilevel:



Example: meta-reinforcement learning

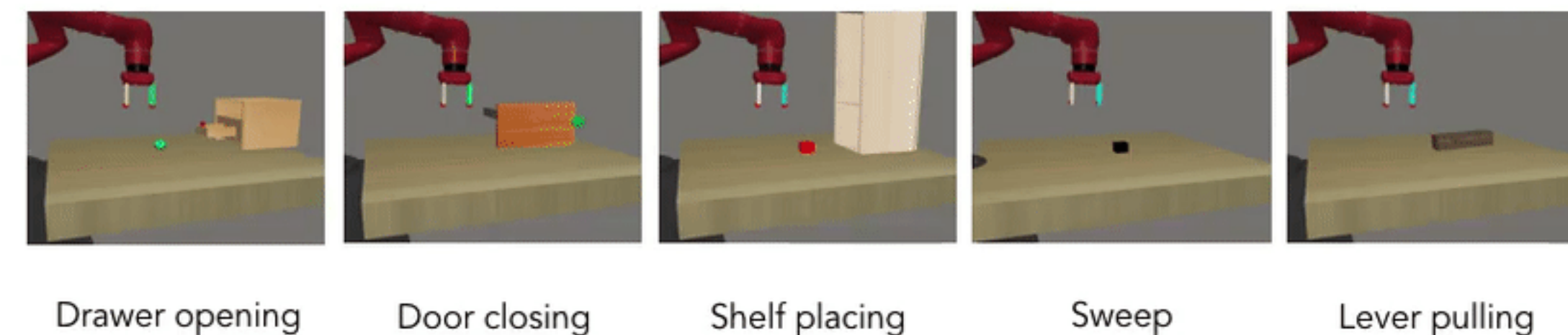
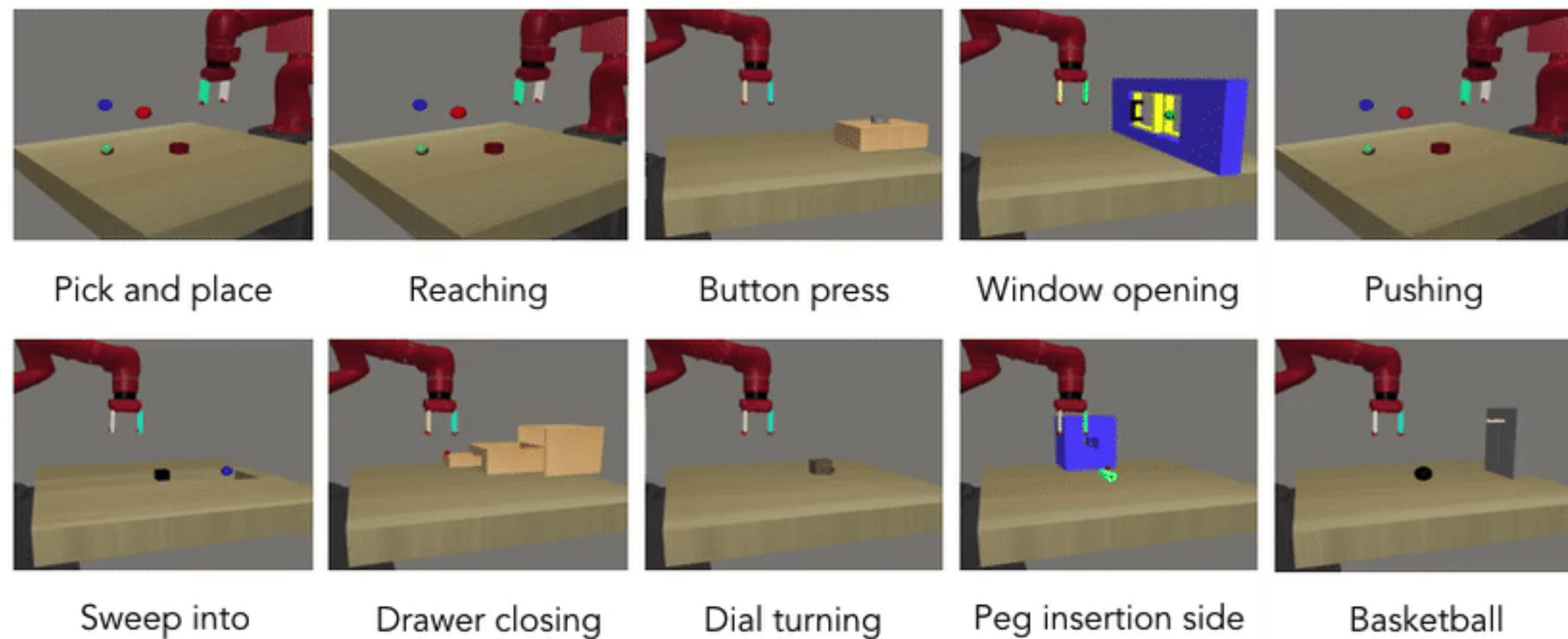
initial state:
randomized object
and goals positions

Train tasks

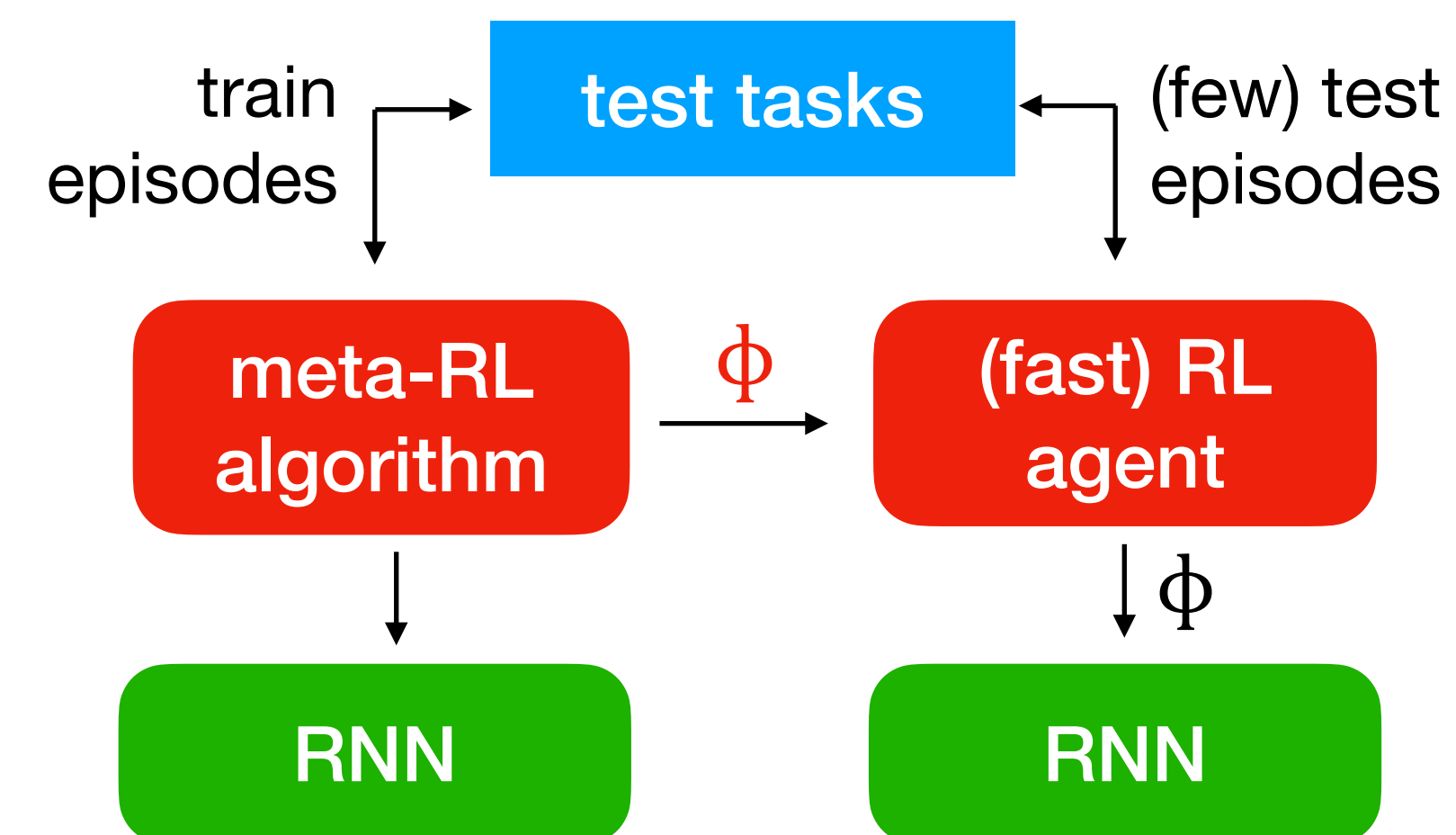
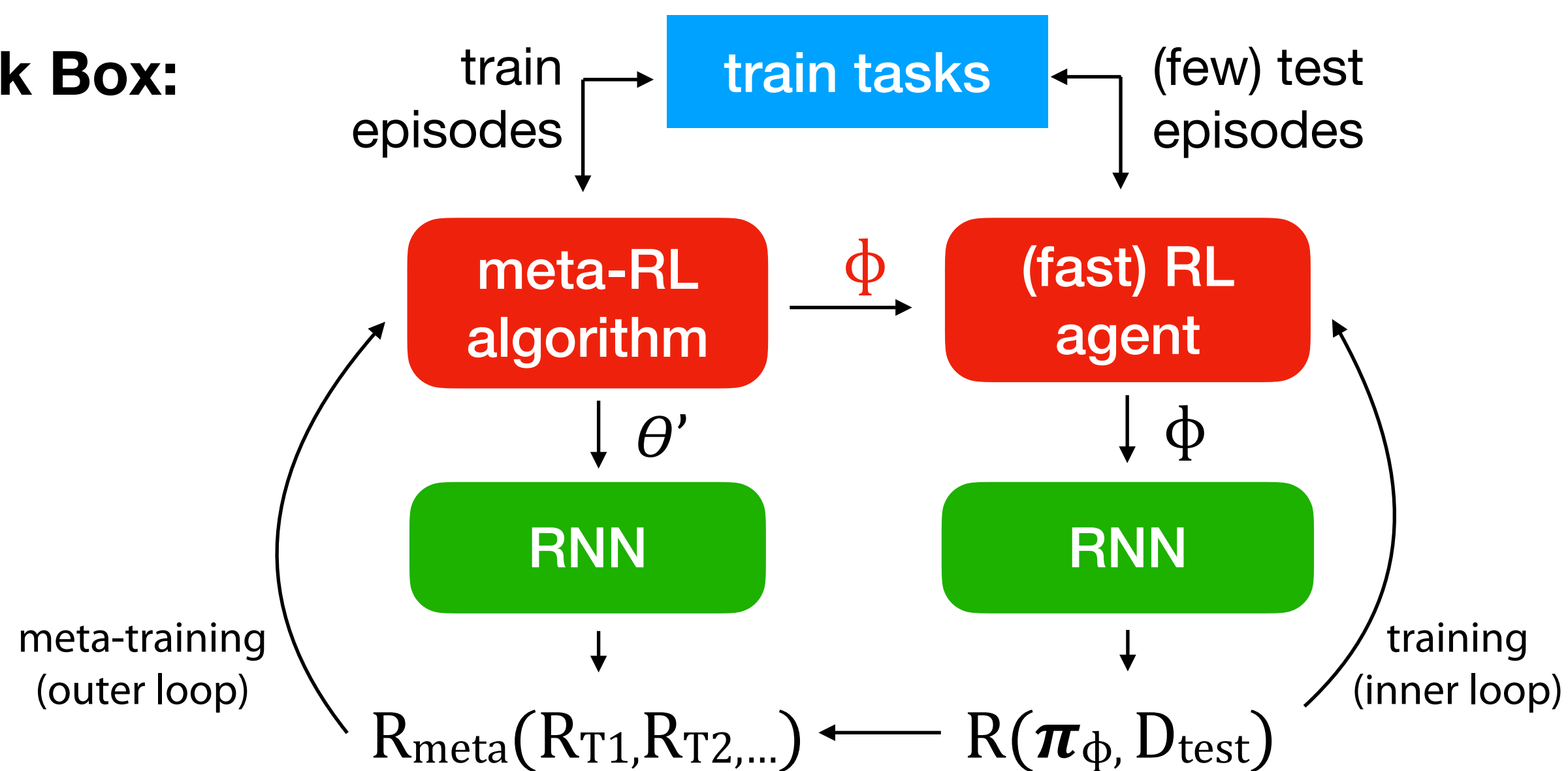
Test tasks

other initial states
or related tasks

ML10



Black Box:



Example: meta-reinforcement learning

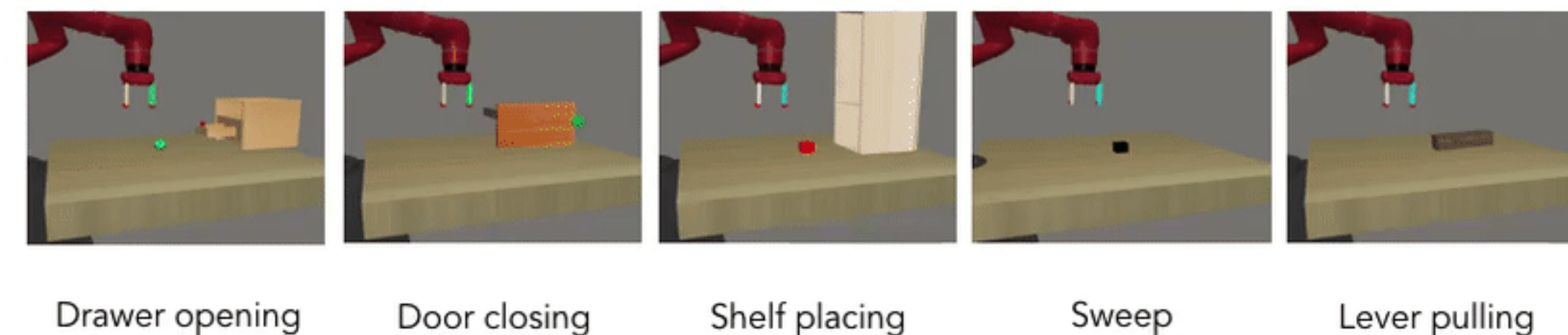
initial state:
randomized object
and goals positions

Train tasks

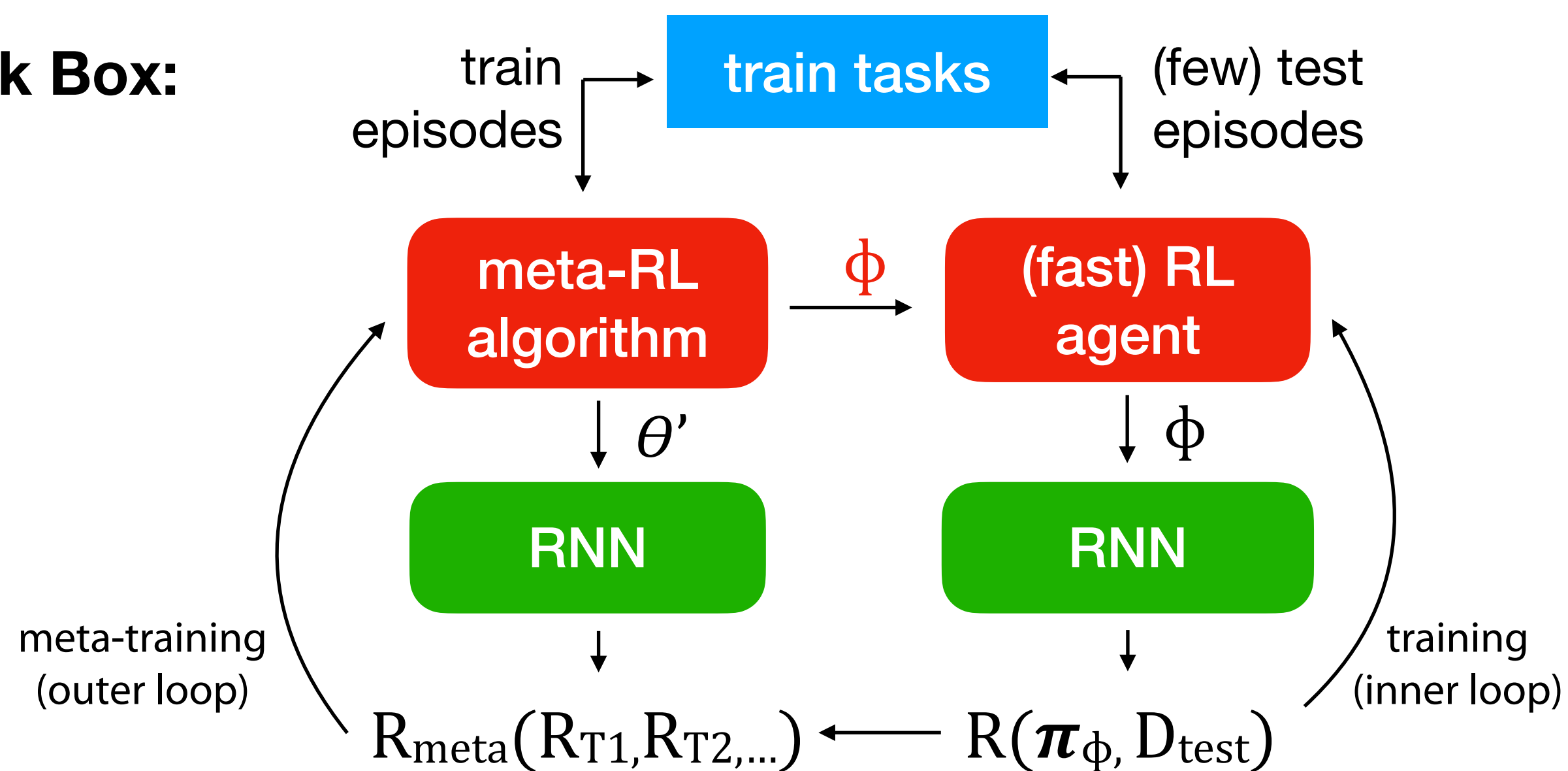
Test tasks

other initial states
or related tasks

ML10

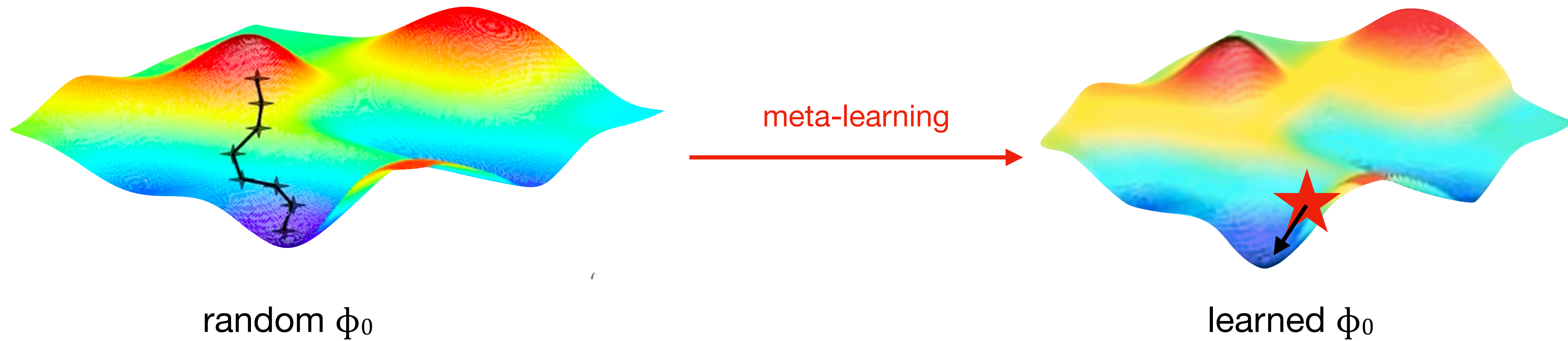
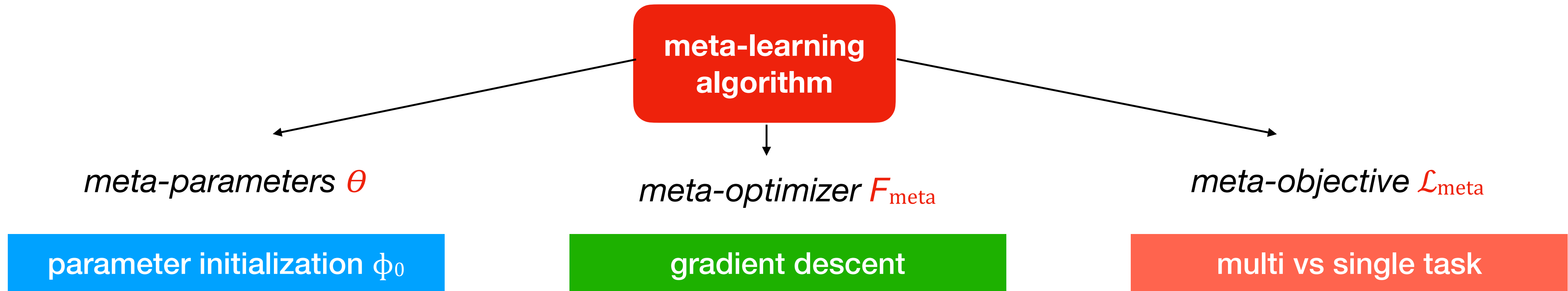


Black Box:



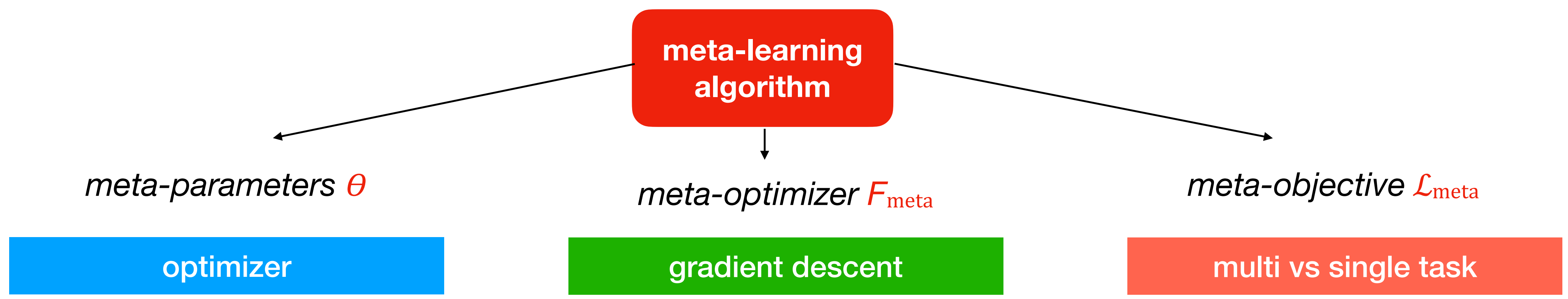
Taxonomy of meta-learning methods

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



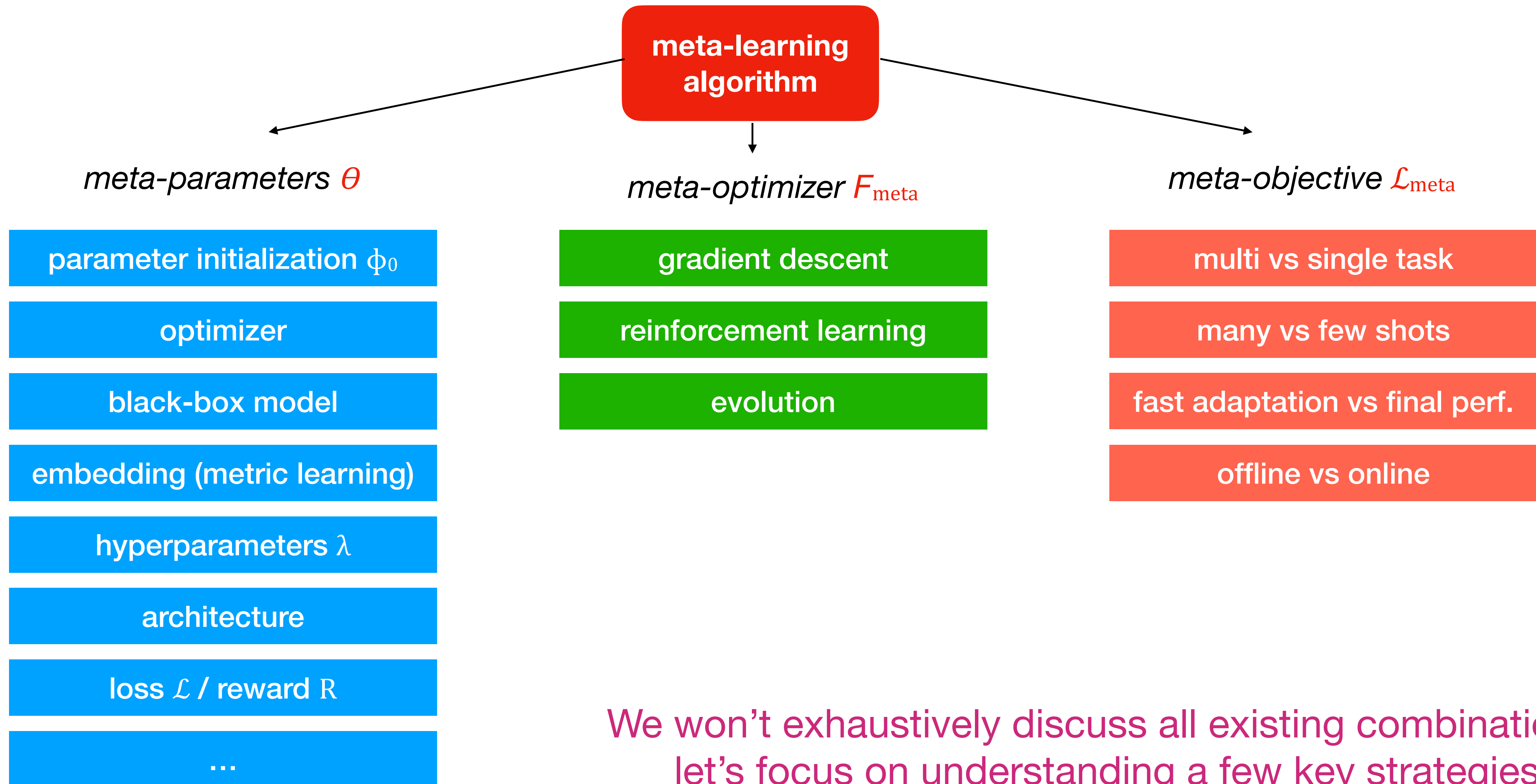
Taxonomy of meta-learning methods

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

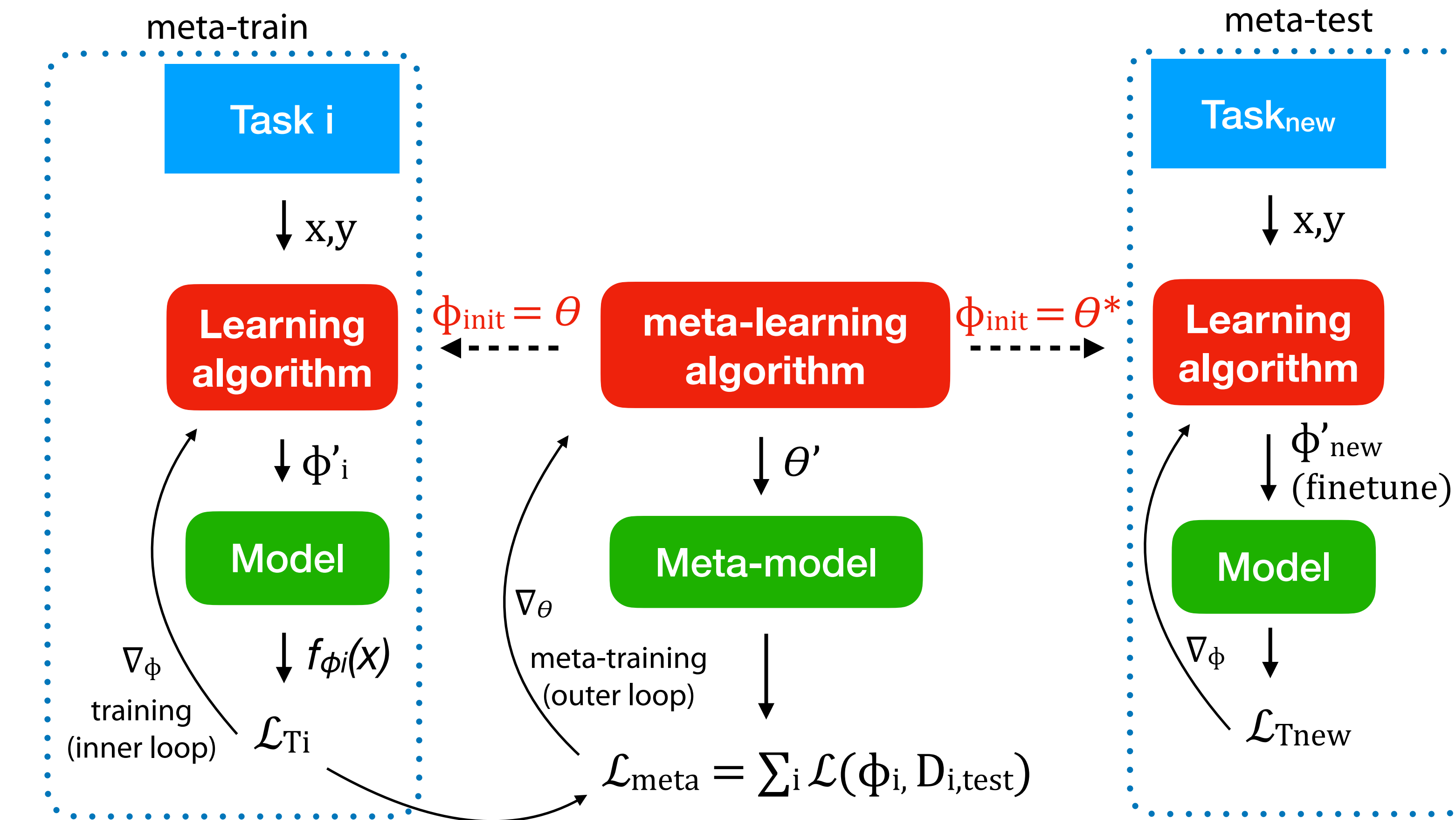


Taxonomy of meta-learning methods

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



Gradient-based methods: learning ϕ_{init}



- θ (prior): model initialization ϕ_{init}
 - learn representation suitable for many tasks (e.g. pretrained CNN)
 - maximize rapid learning
- Each task i yields task-adapted ϕ_i
 - Update algorithm u
- Finetune θ^* on T_{new} (in few steps)

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

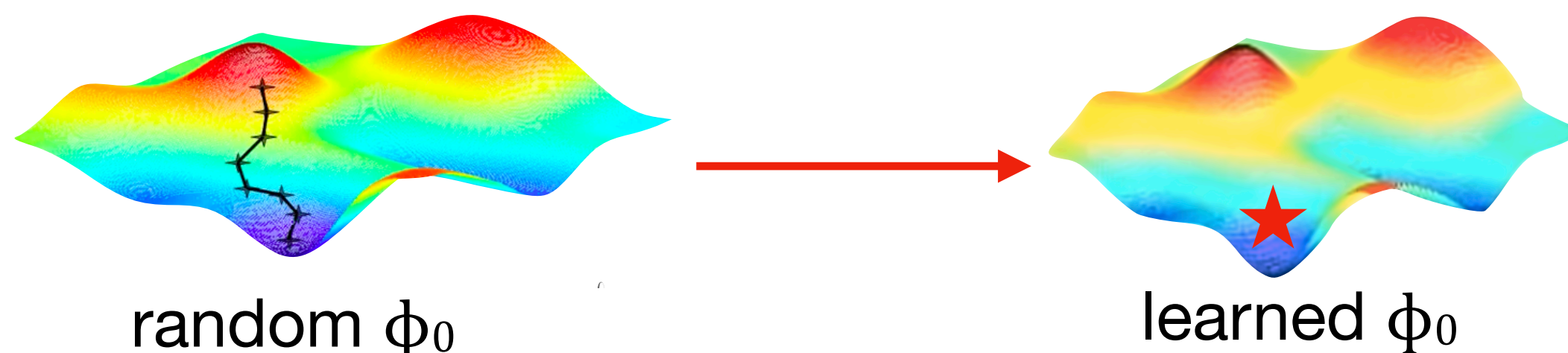
$$\phi'_{\text{new}} = u(\theta^*, D_{\text{new},\text{train}})$$

Can be seen as bilevel optimization:

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

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

$p(T)$



Model agnostic meta-learning (MAML)

Meta-training

- Current initialization θ , model f_θ
- On i tasks, perform k SGD steps to find ϕ_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 θ to minimize sum of per-task losses, repeat

$$\theta \leftarrow \theta - \beta \nabla_\theta \sum_i \mathcal{L}_i(f_{\phi_i}) \quad \alpha, \beta: \text{learning rates}$$

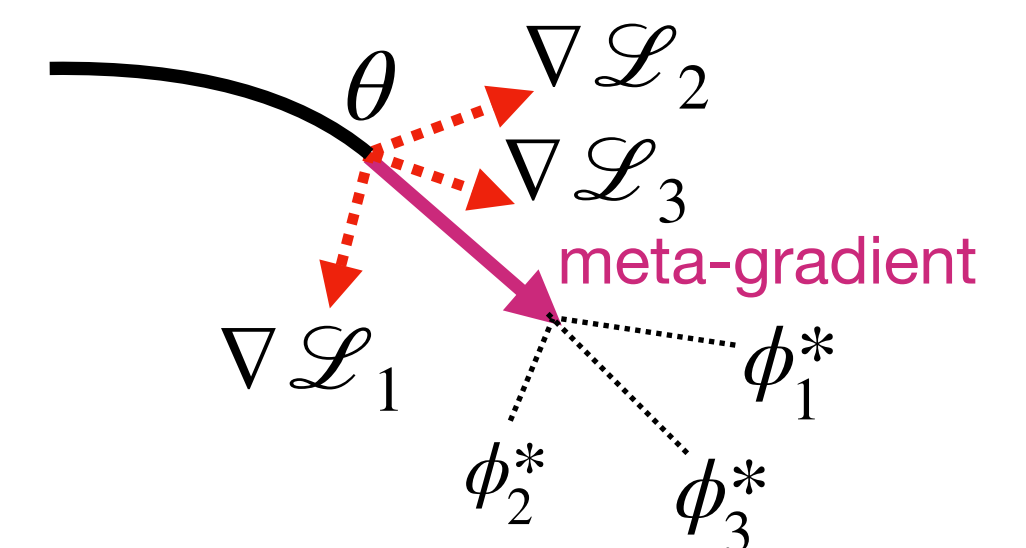
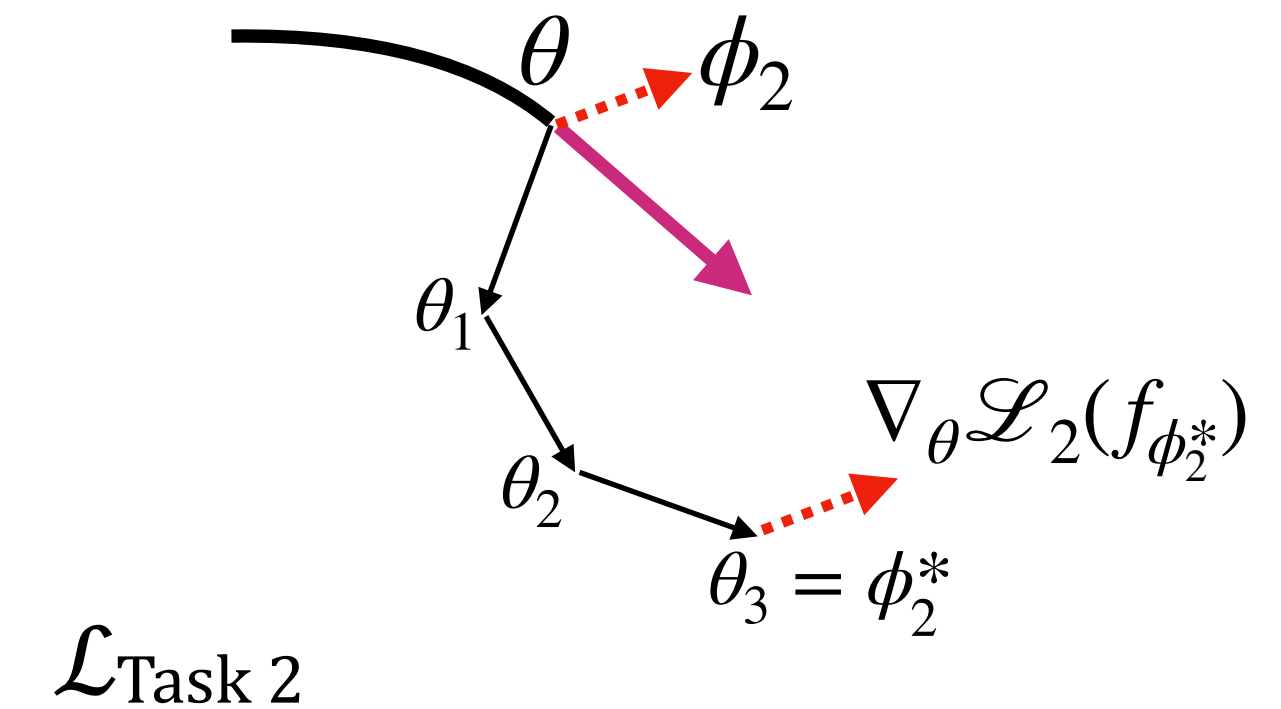
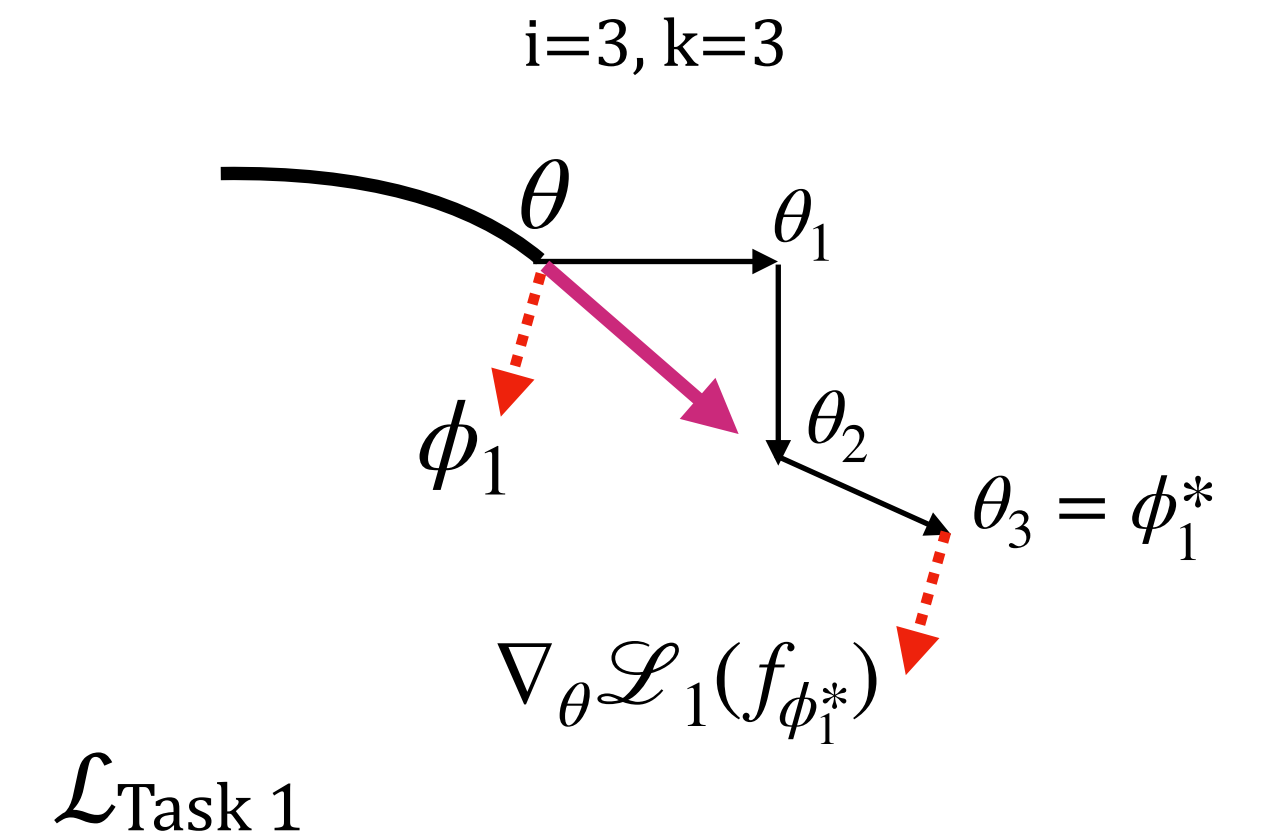
$$\theta \leftarrow \theta - \beta \nabla_\theta \sum_i \mathcal{L}_i(f(\theta - \alpha \nabla_\theta \mathcal{L}_i(f_{\phi_i^*})))$$

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

Meta-testing

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

derivative of test-set loss



Model agnostic meta-learning (MAML)

Meta-training

- Current initialization θ , model f_θ
- On i tasks, perform k SGD steps to find ϕ_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 θ to minimize sum of per-task losses, repeat

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

α, β : learning rates

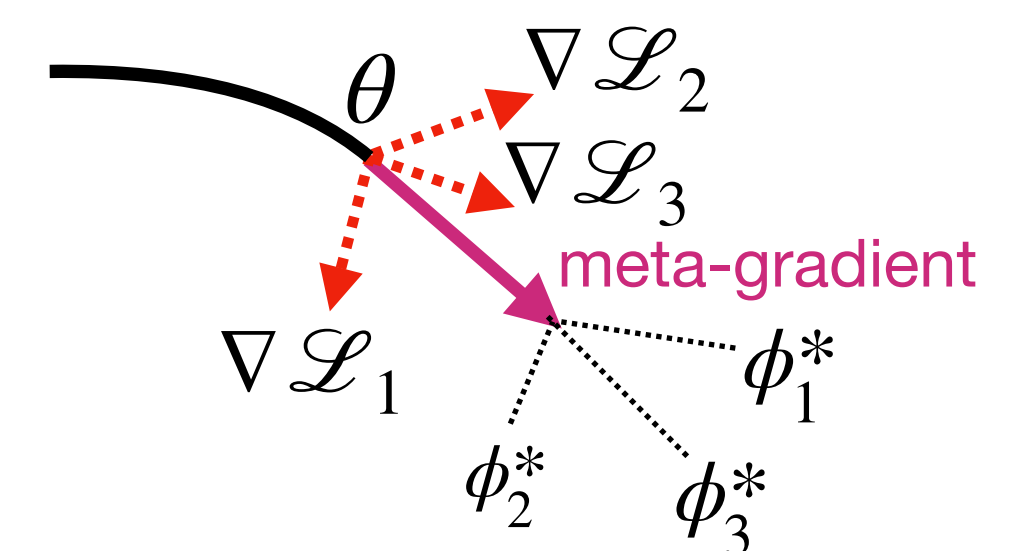
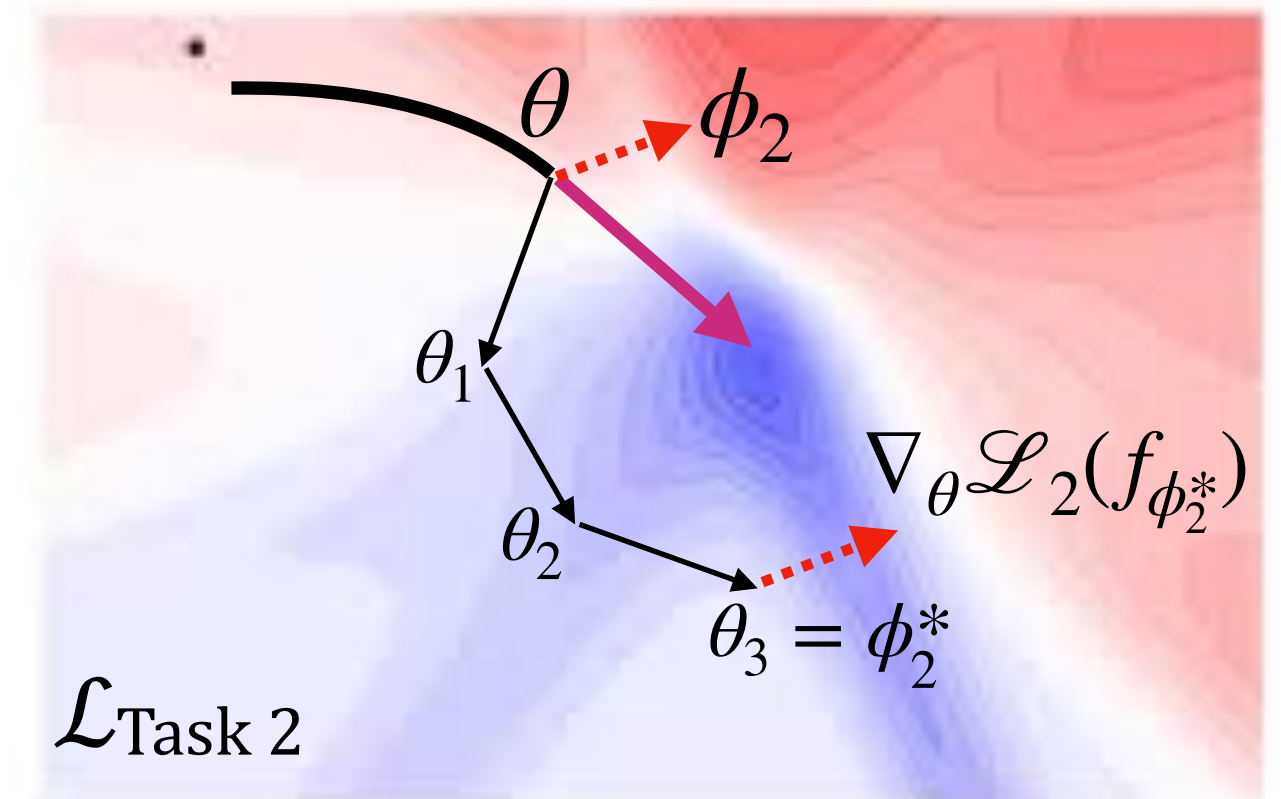
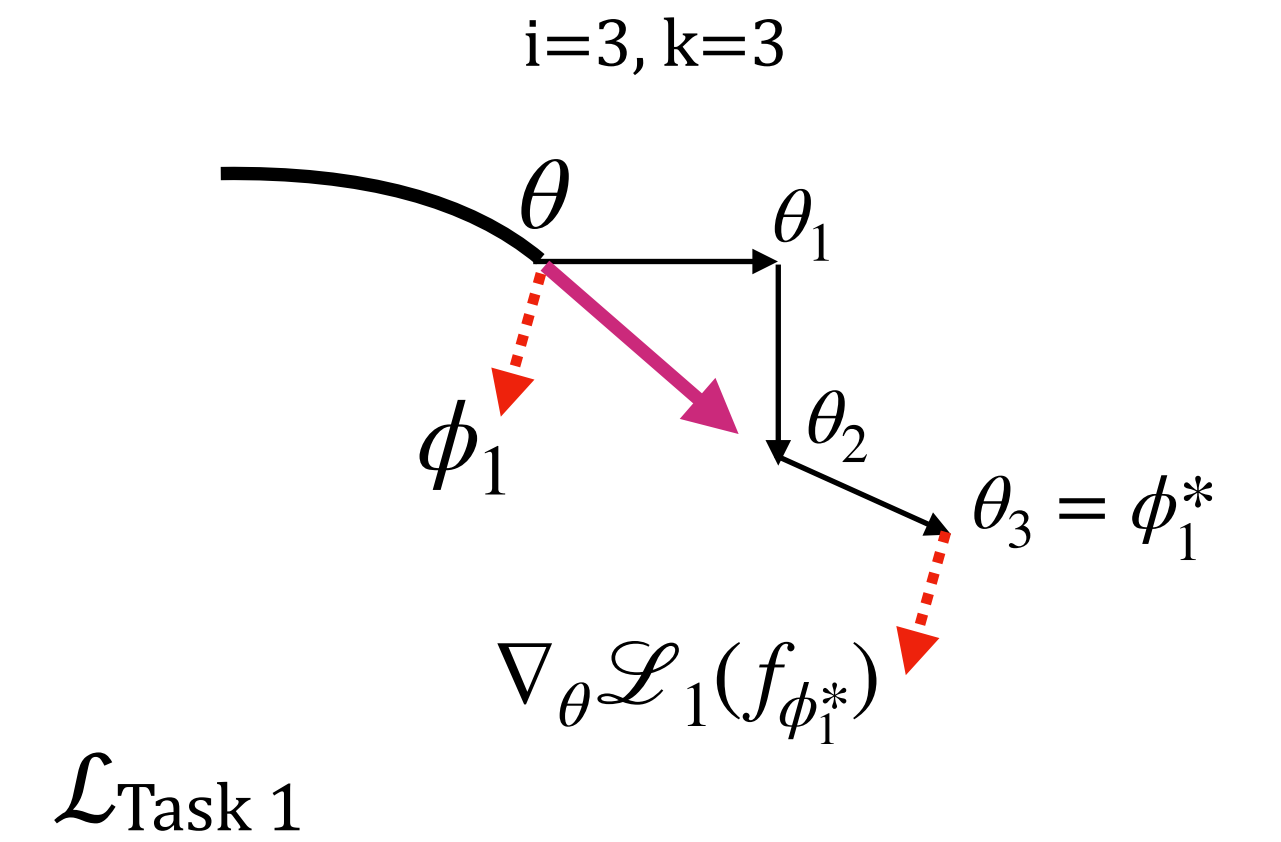
$$\theta \leftarrow \theta - \beta \nabla_\theta \sum_i \mathcal{L}_i(f(\theta - \alpha \nabla_\theta \mathcal{L}_i(f_{\phi_i^*})))$$

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

Meta-testing

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

derivative of test-set loss



Model agnostic meta-learning (MAML)

Meta-training

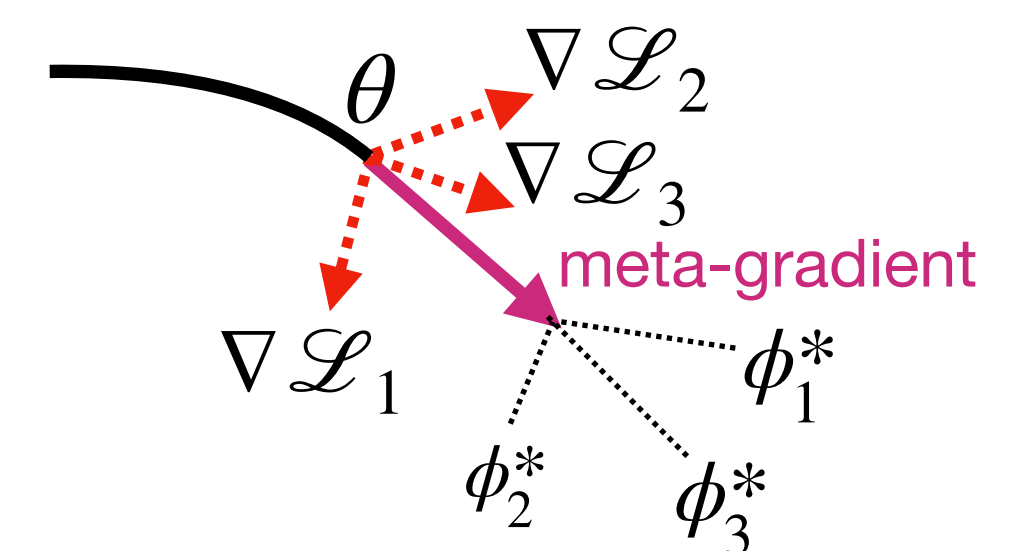
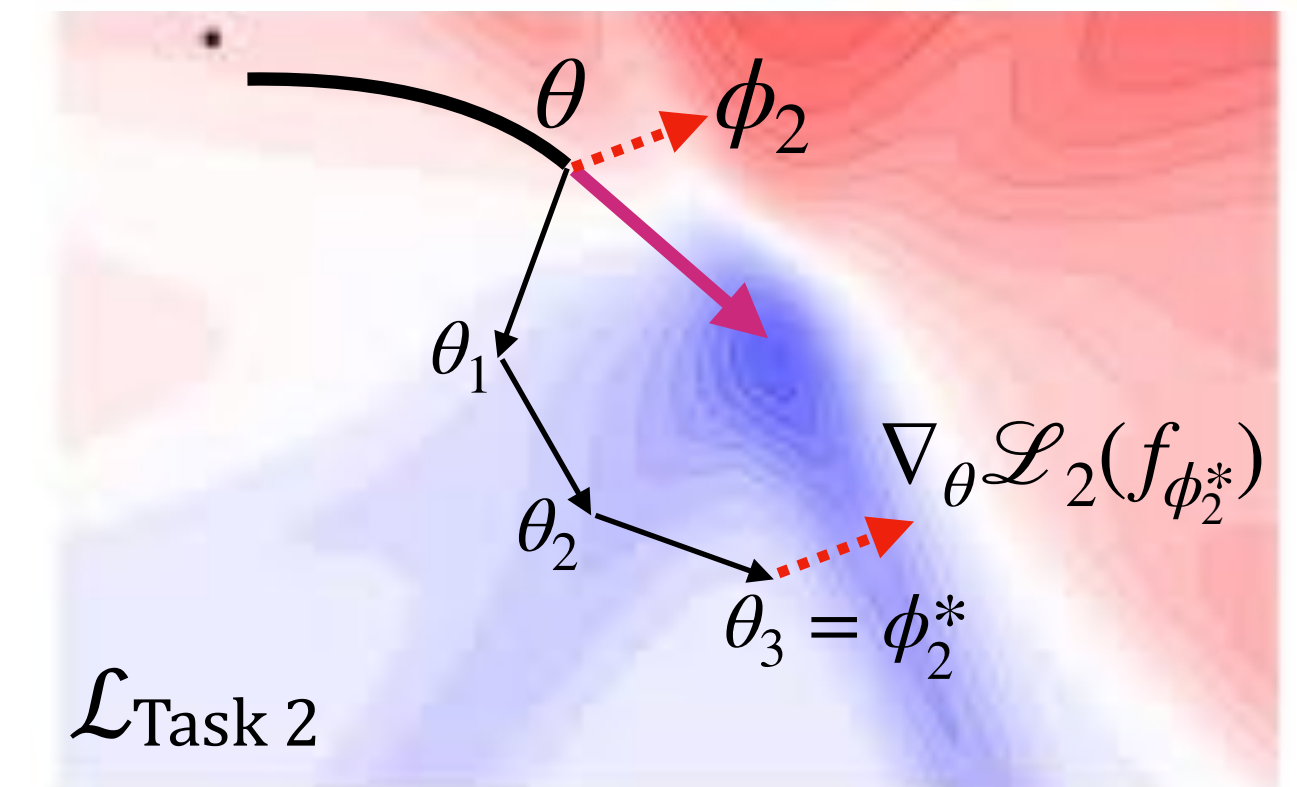
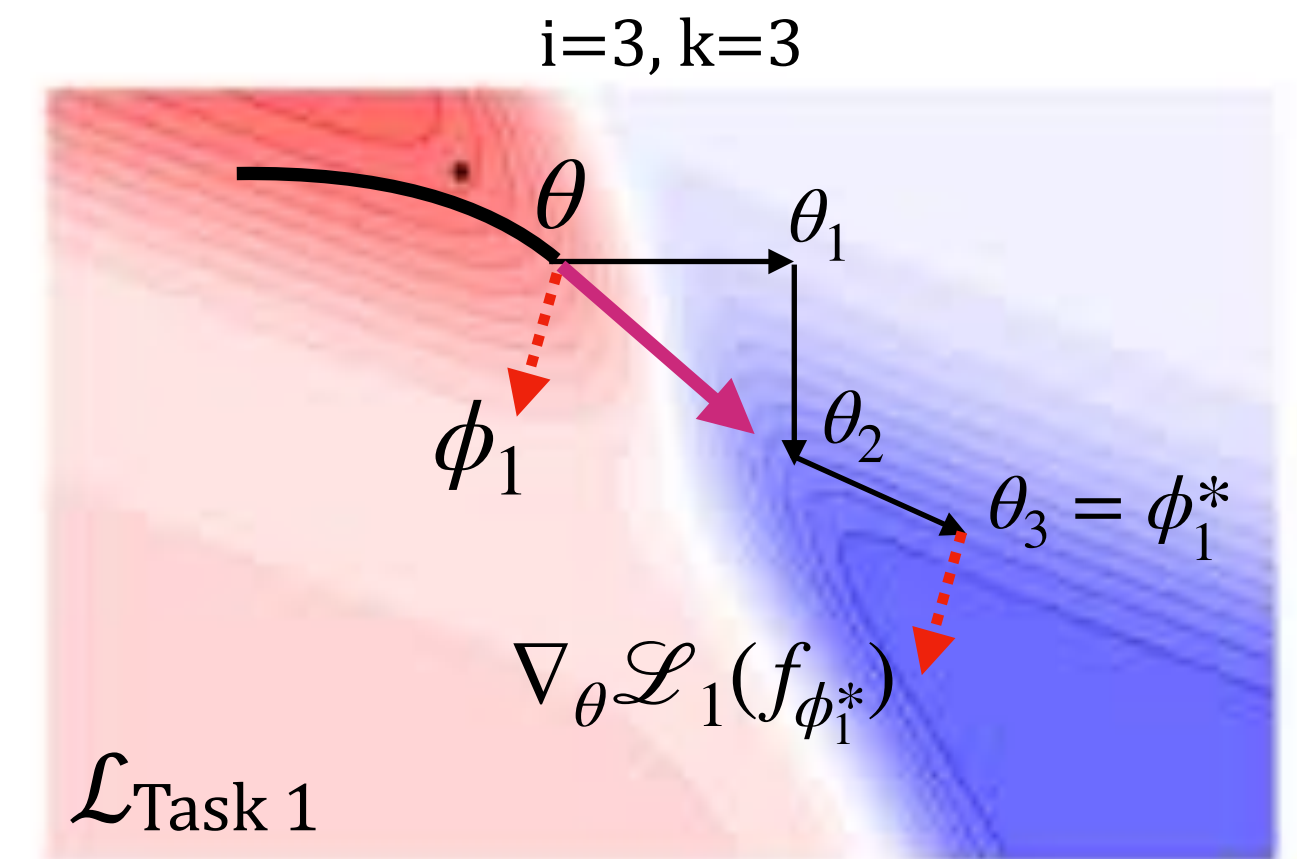
- Current initialization θ , model f_θ
- On i tasks, perform k SGD steps to find ϕ_i^* , then evaluate $\nabla_\theta \mathcal{L}_i(f_{\phi_i^*})$
- Update task-specific parameters: $\phi_i = \theta - \alpha \nabla_\theta \mathcal{L}_i(f_{\phi_i^*})$
- Update θ to minimize sum of per-task losses, repeat

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

α, β : 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 θ affect the gradient at new θ

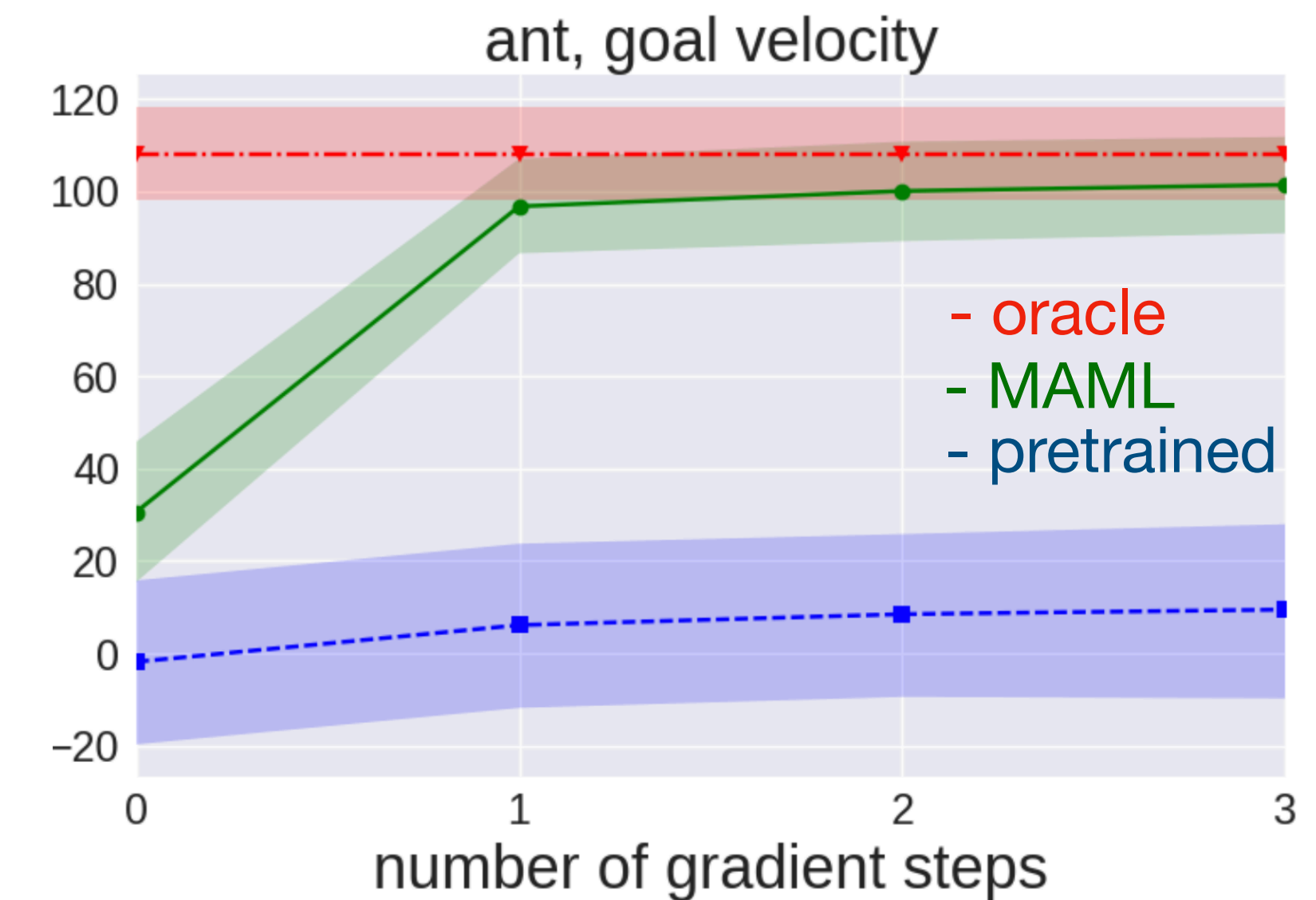
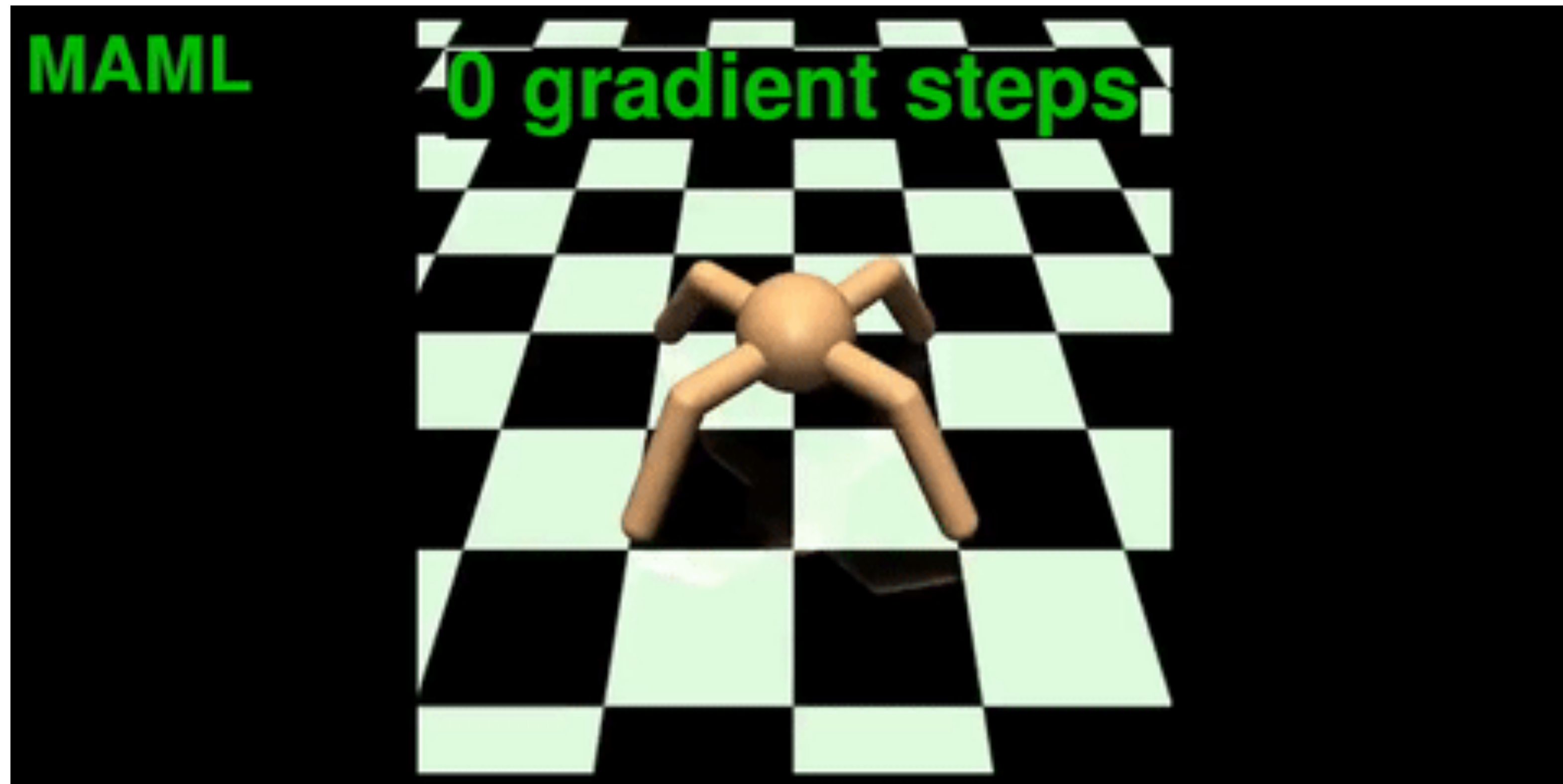


Meta-testing

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

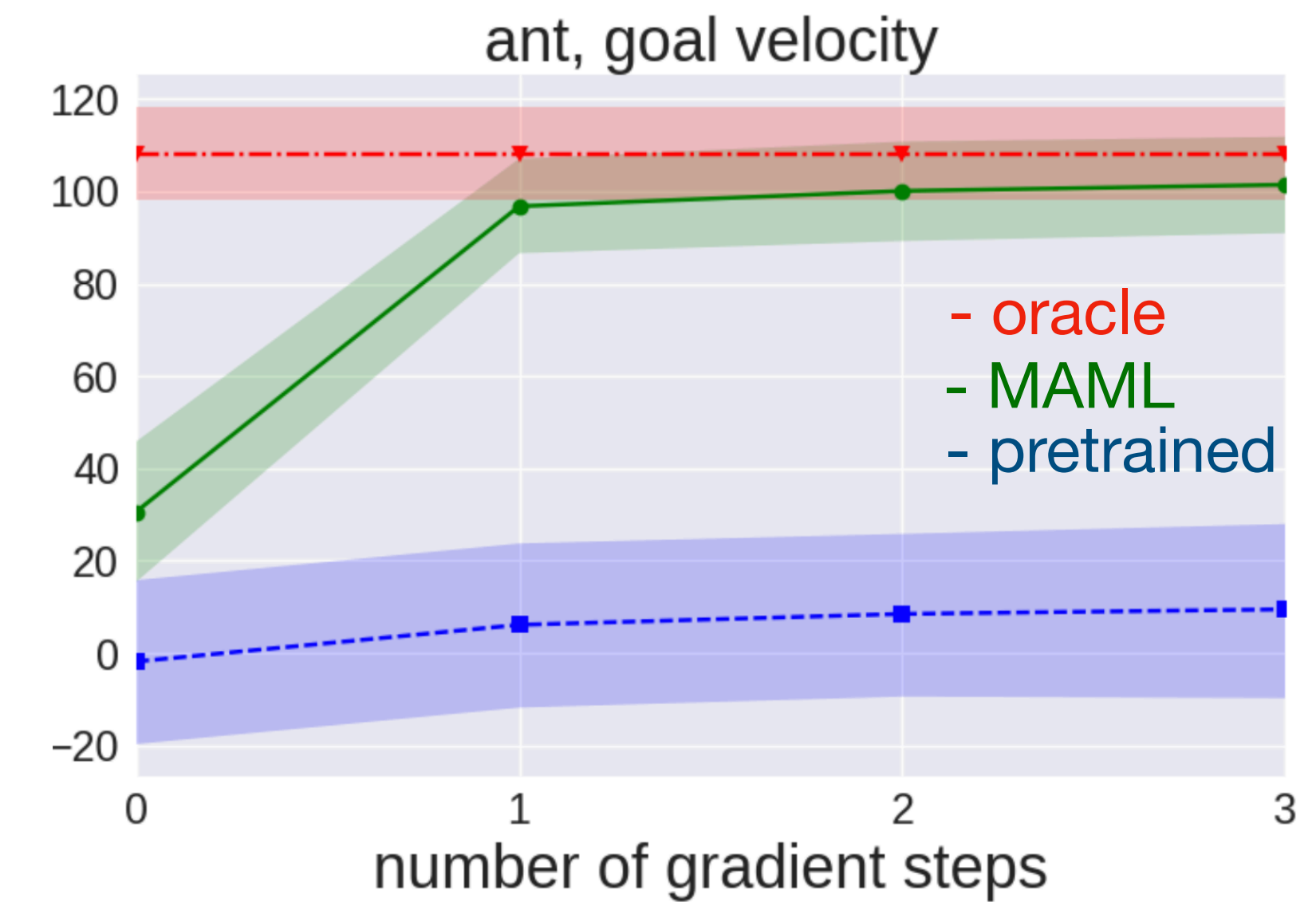
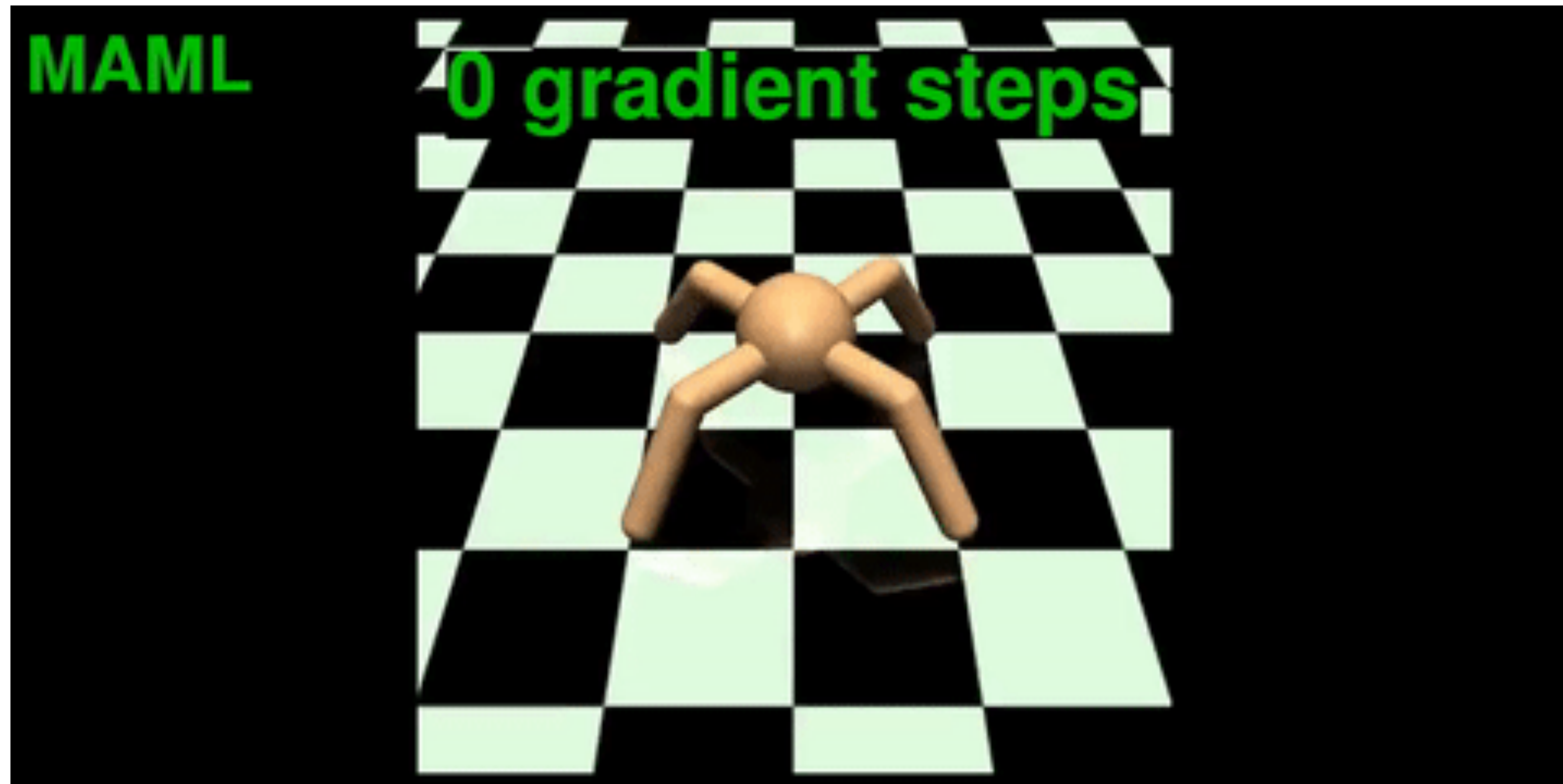
Model agnostic meta-learning (MAML)

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



Model agnostic meta-learning (MAML)

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



Other gradient-based techniques

¹ [Finn et al. 2017](#)

² [Li et al. 2017](#)

³ [Lee et al. 2018](#)

⁴ [Park et al. 2019](#)

⁵ [Flennerhag et al. 2019](#)

⁶ [Antoniou et al. 2019](#)

⁷ [Finn et al. 2019](#)

- Changing update rule yield different variants:

- MAML ^{1,6} $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(f_{\phi_i})$

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

- MetaSGD ² $\phi_i = \theta - \alpha \text{diag}(w) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$

- Tnet ³ $\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*}, w)$

- Meta curvature ⁴ $\phi_i = \theta - \alpha B(\theta, w) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$

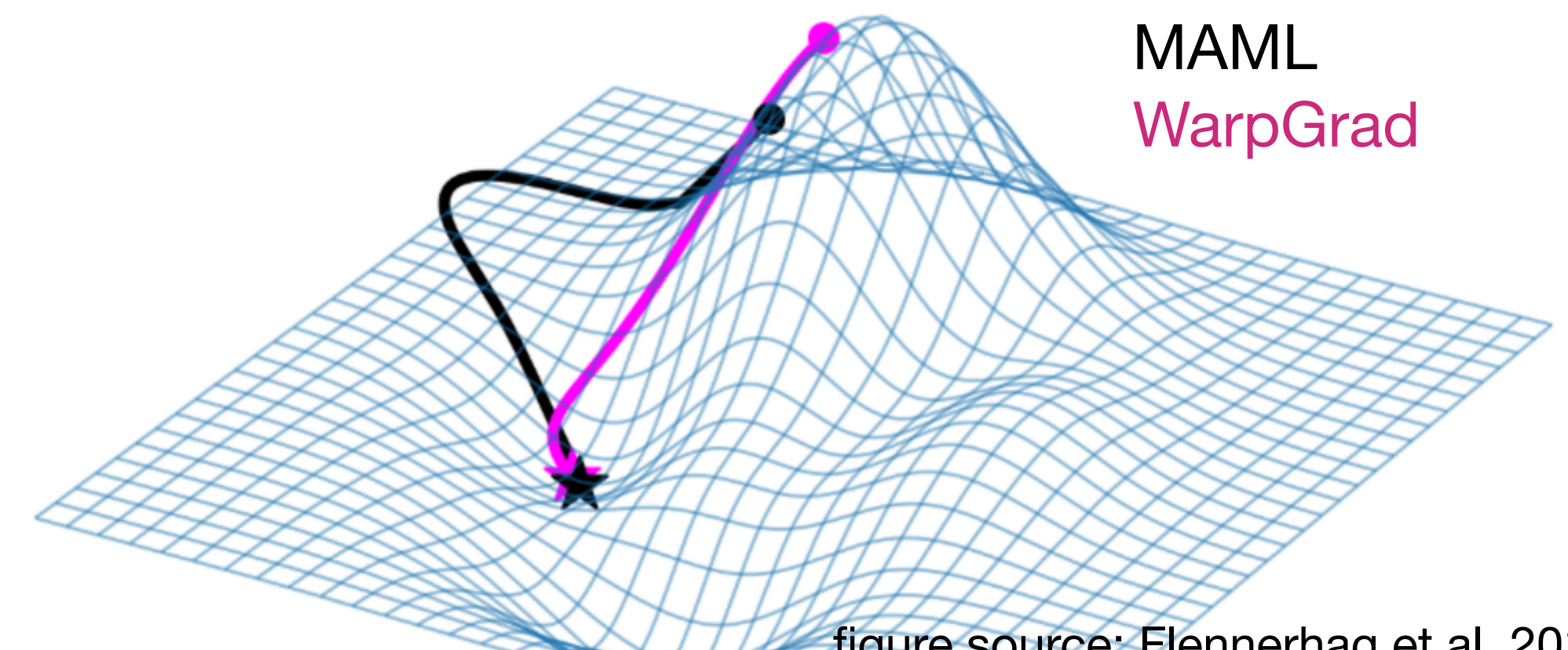
- WarpGrad ⁵ $\phi_i = \theta - \alpha P(\theta, \phi_i) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$

w: weight per parameter

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

- Online MAML (Follow the Meta Leader) ⁷

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



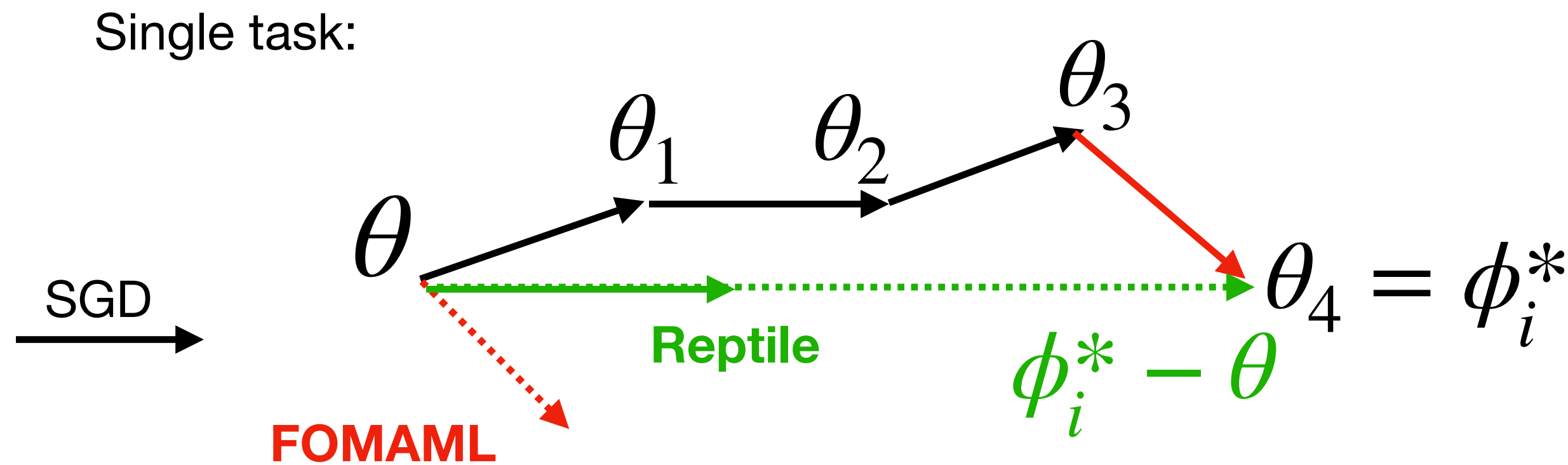
Scalability

¹ Finn et al. 2017

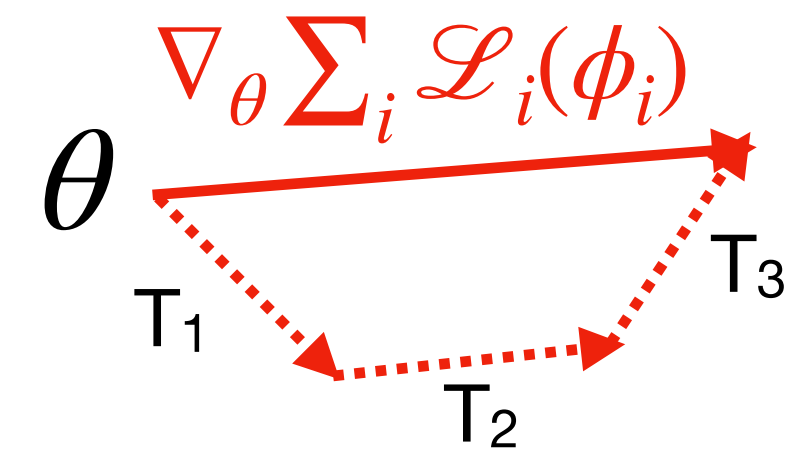
² Nichol et al. 2018

³ Rajeswaran et al. 2019

- Backpropagating derivatives of ϕ_i wrt θ is compute + memory intensive (for large models)
- First order approximations of MAML:
 - First order MAML¹ uses only the last inner gradient update: $\theta \leftarrow \theta - \beta \sum_i \mathcal{L}_i(f_{\phi_i^*})$
 - Reptile²: iterate over tasks, update θ in direction of ϕ_i^* : $\theta \leftarrow \theta - \beta (\phi_i^* - \theta)$

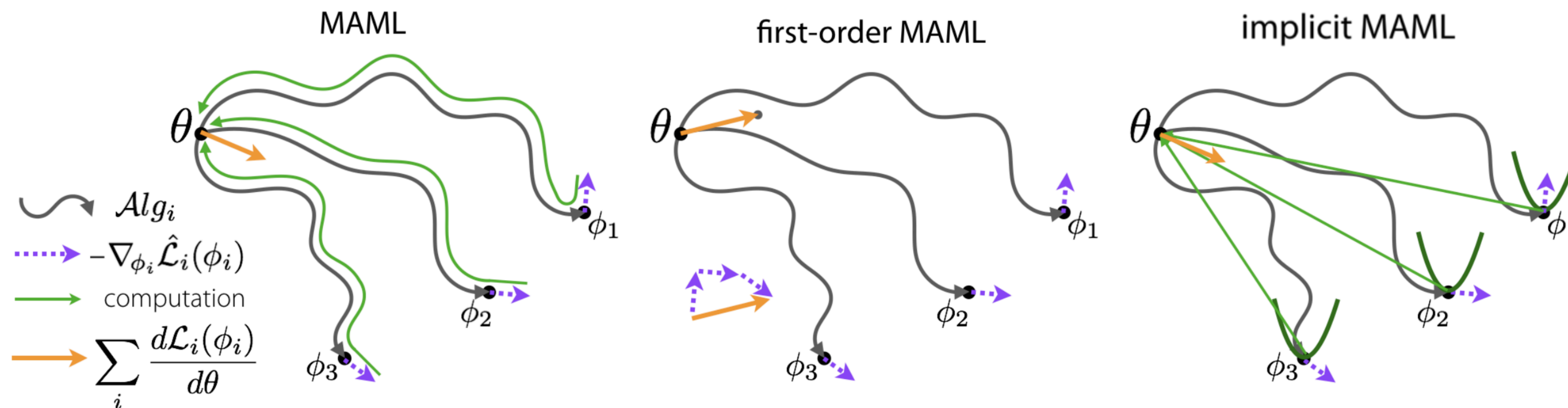


FOMAML meta-gradient:



Scalability (2)

- Implicit MAML¹: uses an approximate gradient
 - Compute derivative of ϕ_i^* wrt θ
 - ϕ_i^* could be anywhere. Hence, add penalty: $\|\phi_i^* - \theta\|^2$
 - Accurate if we stay close to θ
 - $\mathcal{L}_i(\phi_i) + \lambda\|\phi_i^* - \theta\|^2$ has closed form solution



Generalizability

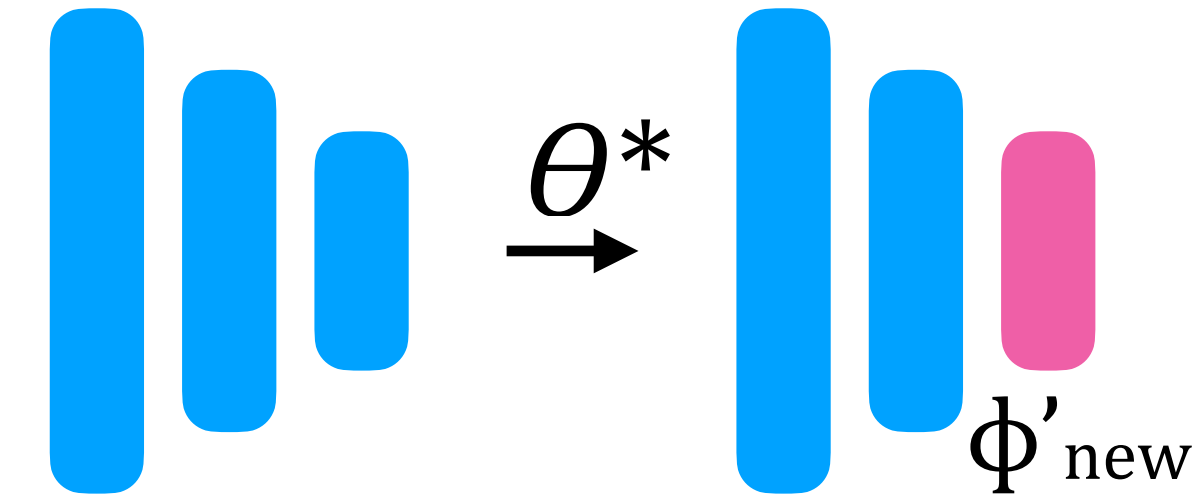
¹ Finn et al. 2018

² Raghu et al. 2020

³ Tian et al. 2020

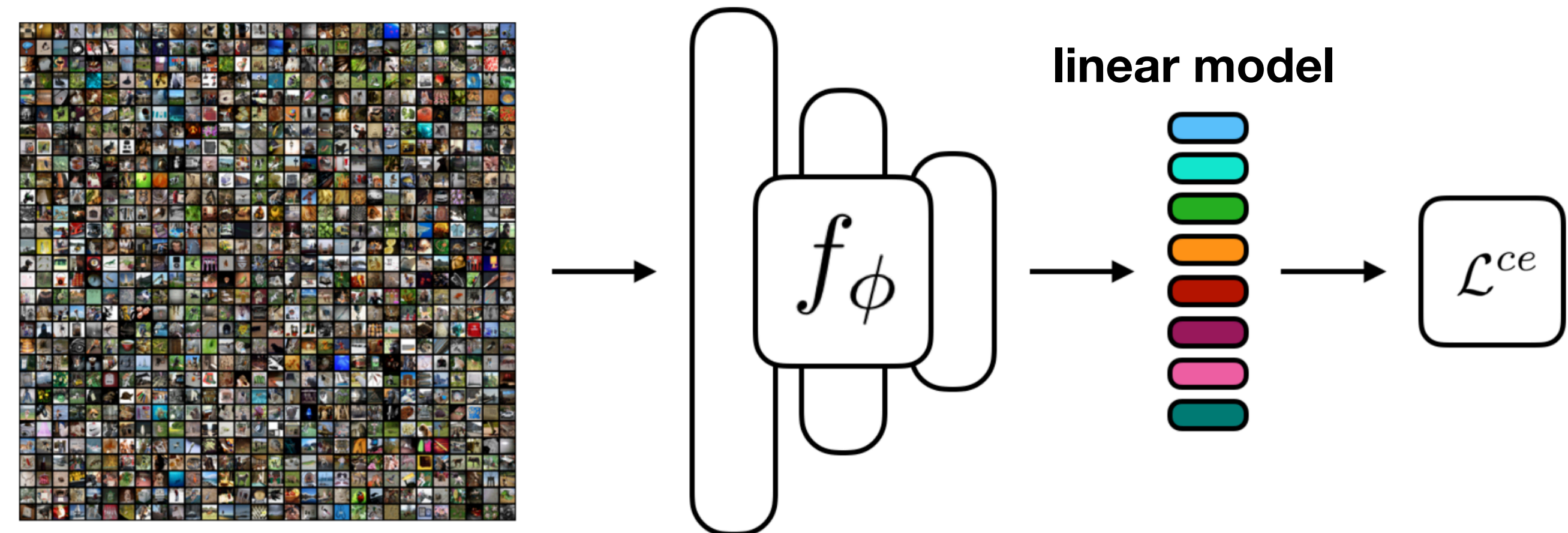
⁴ Stadie et al. 2019

- MAML is more resilient to overfitting than many other meta-learning techniques ¹
- Effectiveness seems mainly due to feature reuse (finding θ) ²
- Fine-tuning *only the last layer* equally good
- On few-shot learning benchmarks, a good embedding outperforms most meta-learning ³



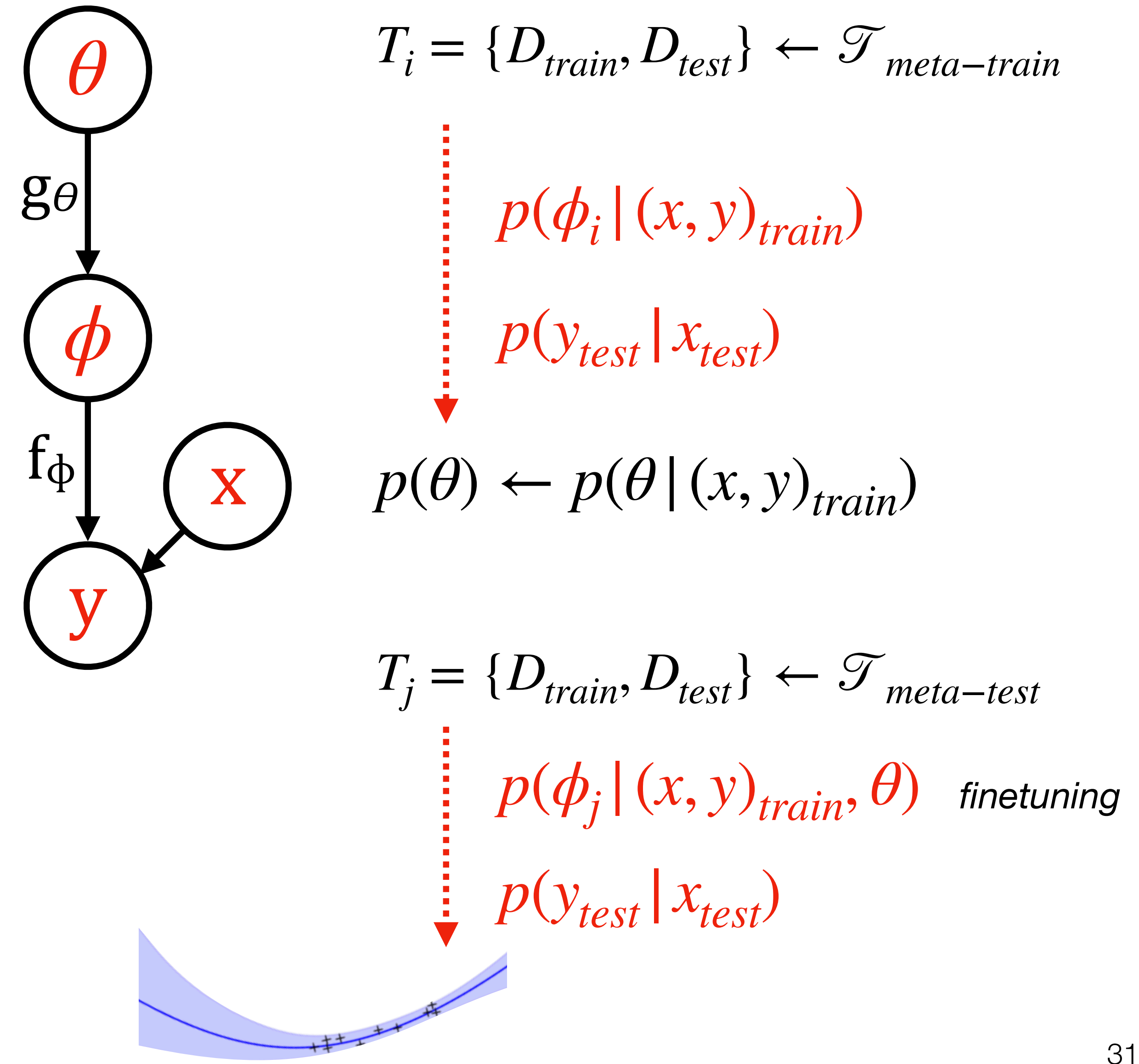
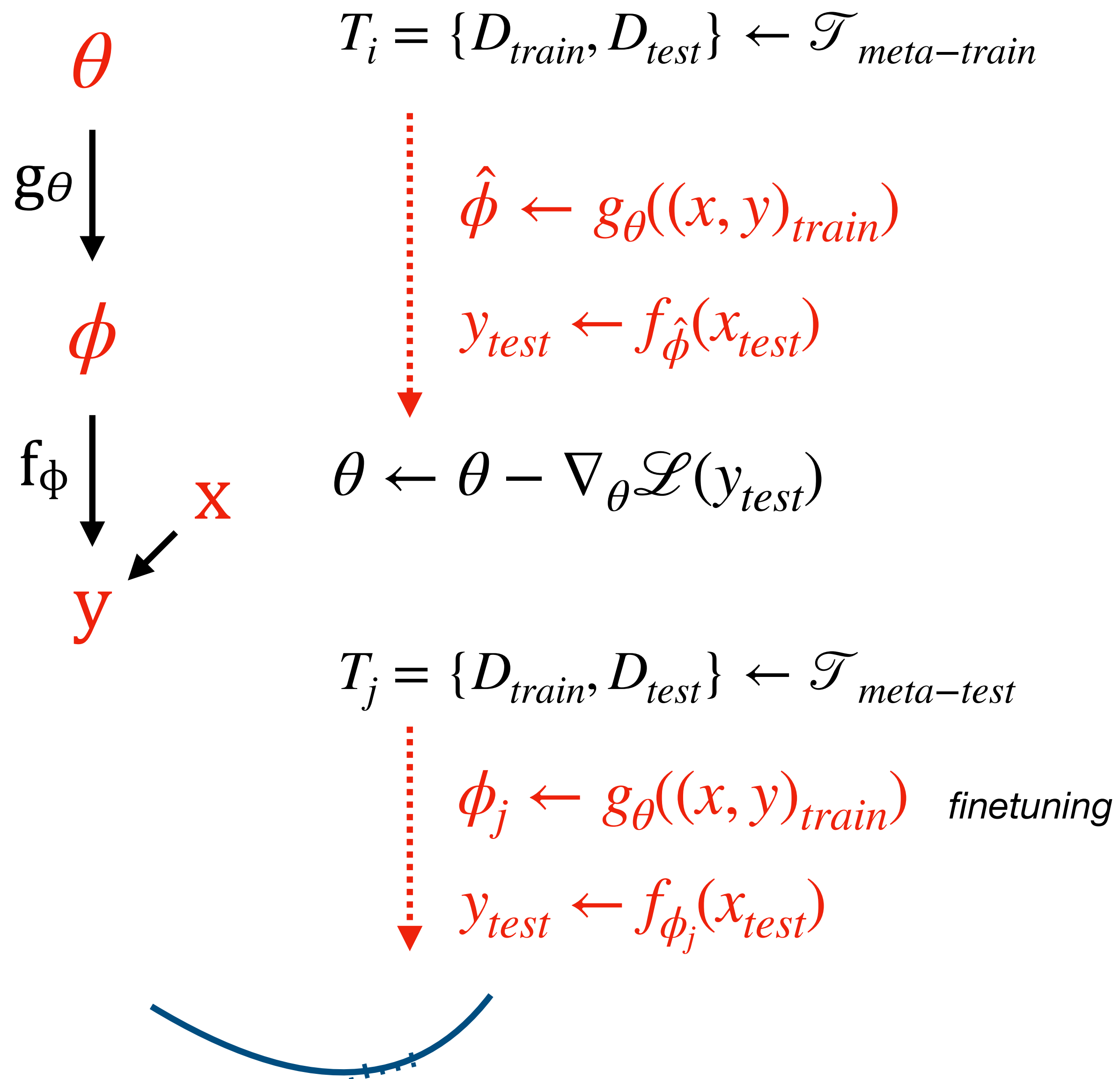
- Learn representation on entire meta-dataset (merged into single task)
- Train a linear classifier on embedded few-shot D_{train} , predict D_{test}

- For meta-RL, also learn how to *explore* (how to sample new environments)
- E-MAML: add exploration to meta-objective (allows longer term goals) ⁴

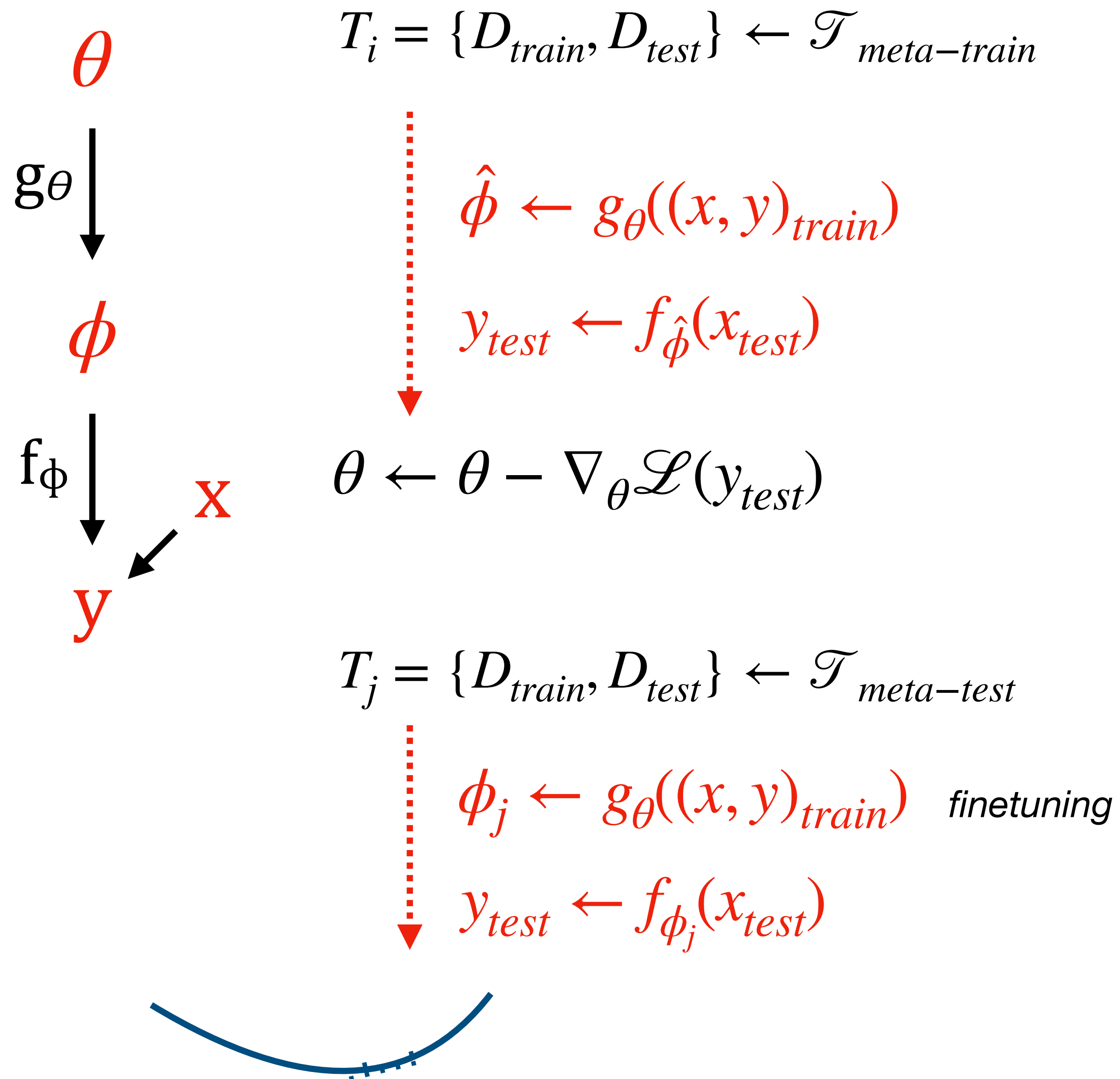


Bayesian meta-learning

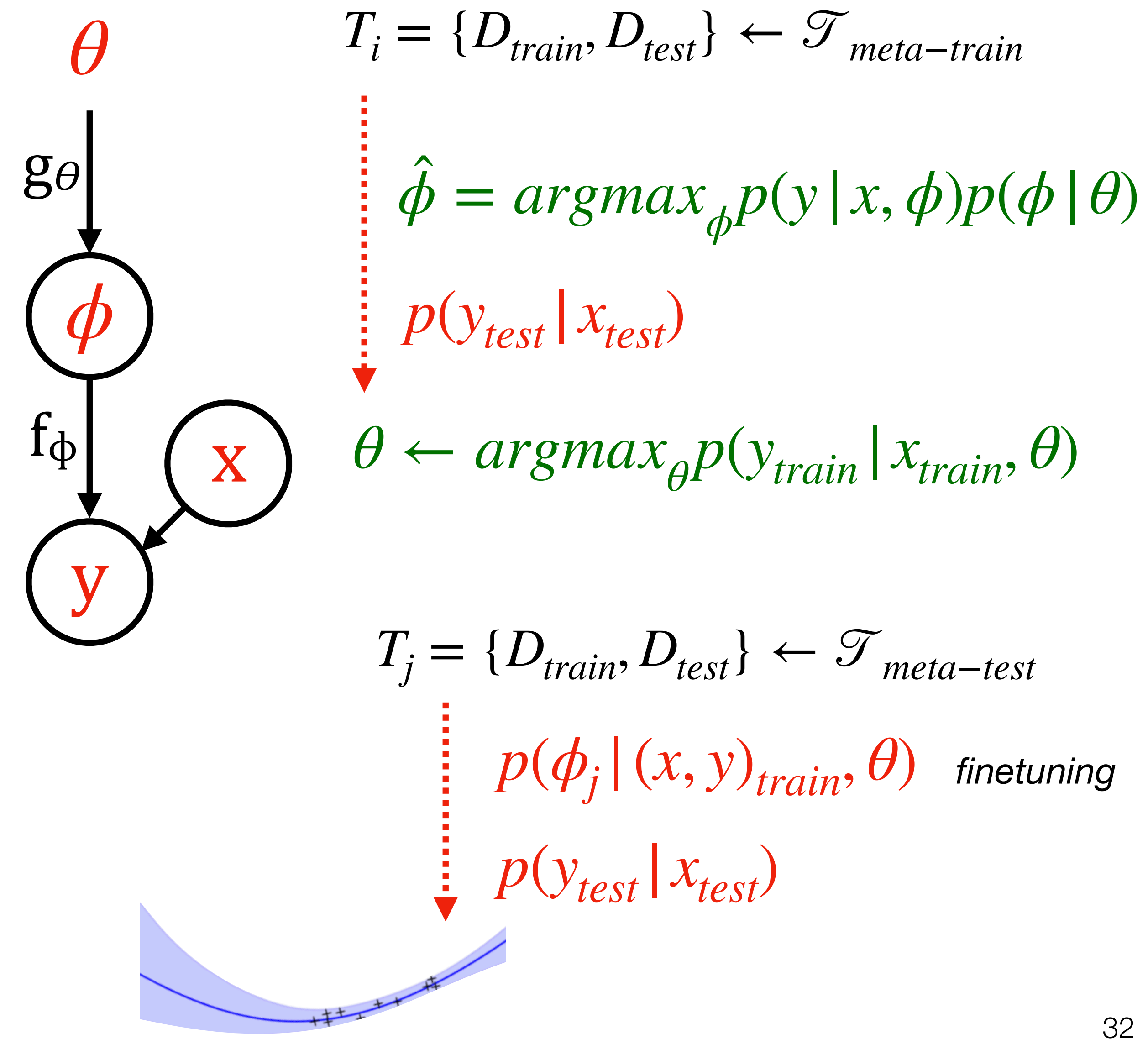
Can meta-learning reason about **uncertainty** in the task distribution?



Bayesian meta-learning



Empirical Bayes: fit point estimate on θ

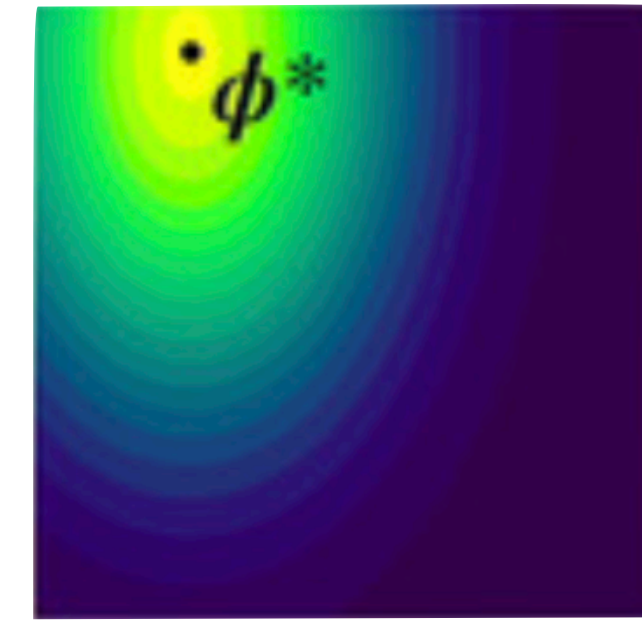


Bayesian meta-learning

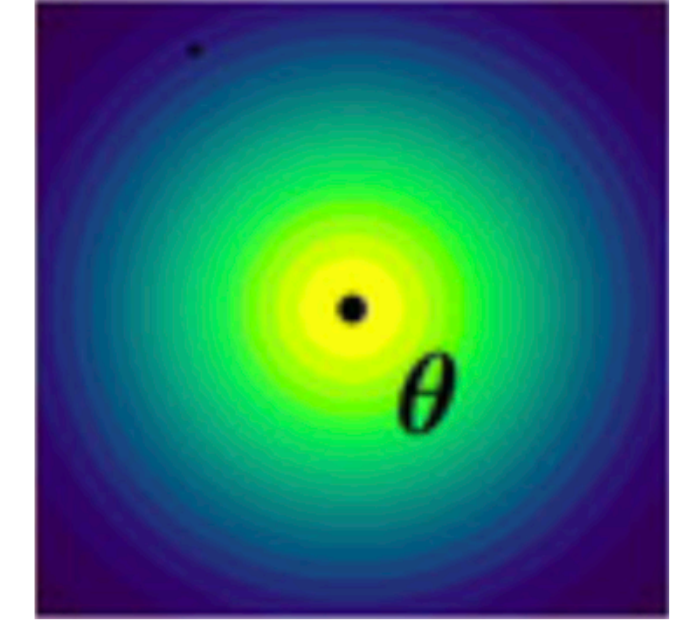
ϕ can be factored as the product of a likelihood and a prior

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

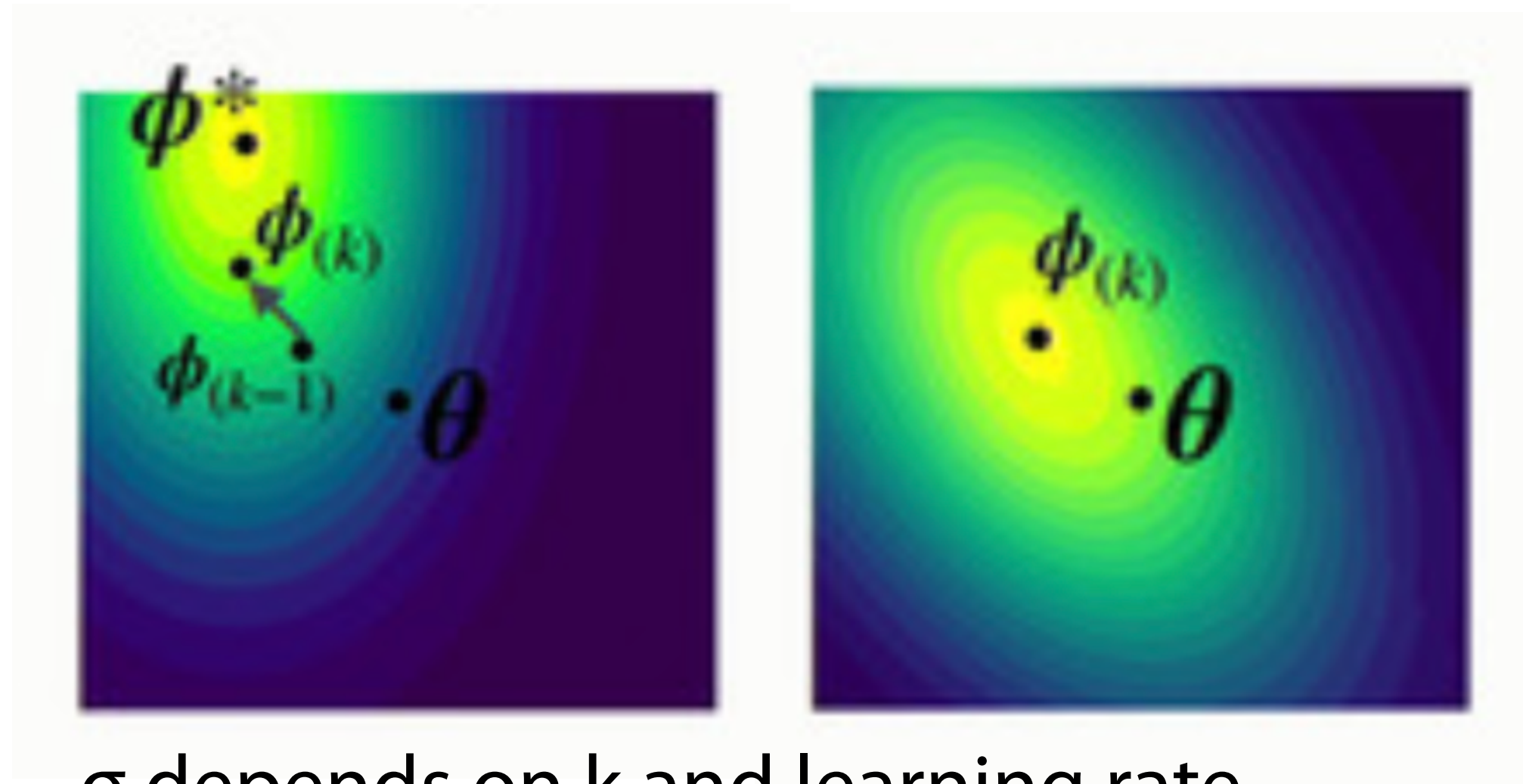
$p(y | x, \phi)$



$p(\phi | \theta)$



When using θ as the initialization of gradient descent, early stopping is approximately equivalent to placing a prior over the task specific parameters ¹



σ depends on k and learning rate

MAML is doing the same

- estimates ϕ_i^* with k steps of gradient descent
- update prior θ to makes it easier to optimize ϕ_i

Fully Bayesian meta-learning

¹ Finn et al. 2019

² Kim et al. 2018

³ Ravi et al. 2019

⁴ Grant et al. 2018

⁵ Edwards et al. 2017

⁶ Garnelo et al. 2018

⁷ Jerfel et al. 2019

- Alternatively, use approximation methods to represent uncertainty over ϕ_i

- Sampling technique + variational inference: PLATIPUS ¹, BMAML ², ABML ³

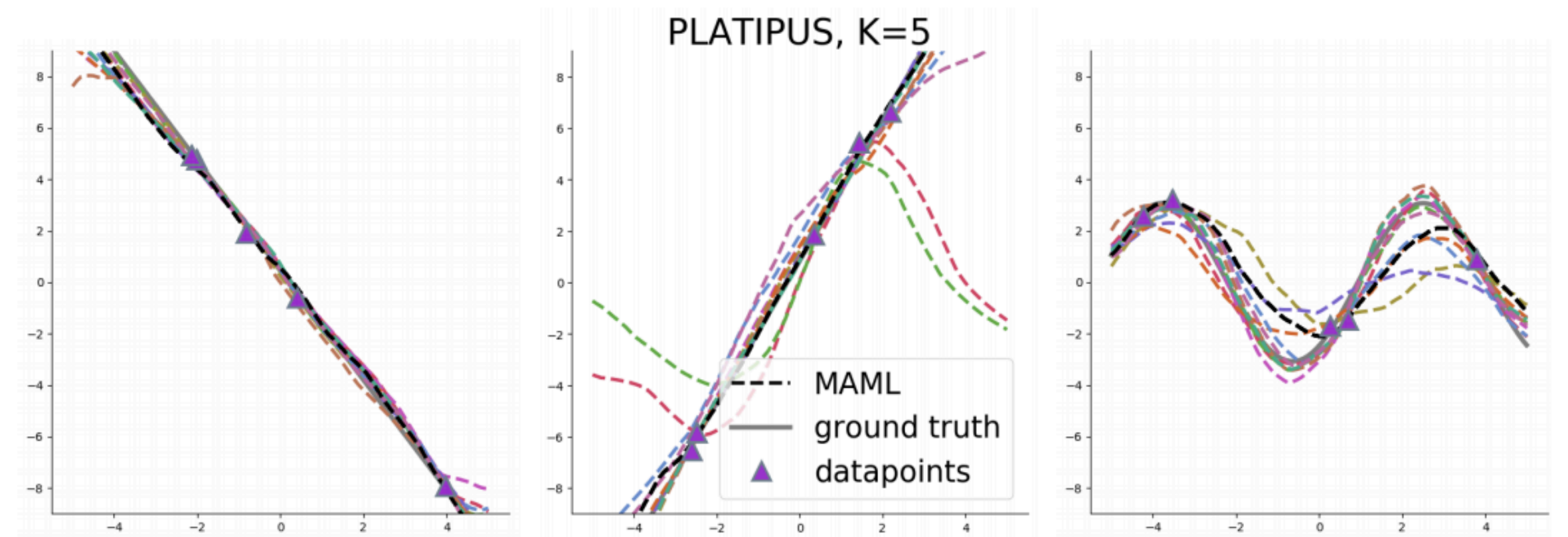
- Laplace approximation: LLAMA ⁴

- Variational approximation of posterior:

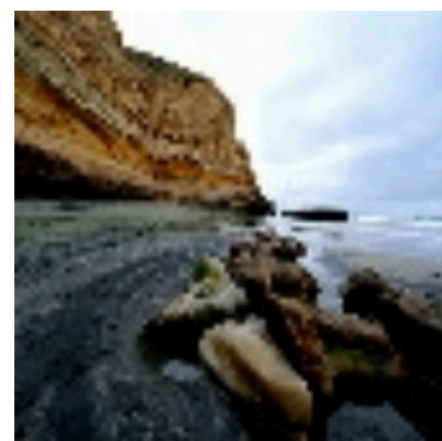
- Neural Statistician ⁵, Neural Processes ⁶

- What if our tasks are not IID?

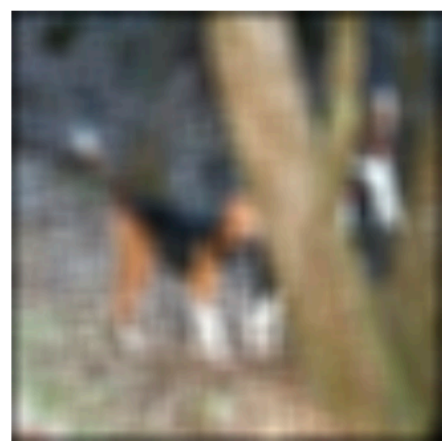
- Impose additional structure, e.g. with task-specific variable z ⁷



mini-ImageNet with filters for non-homogeneous tasks:



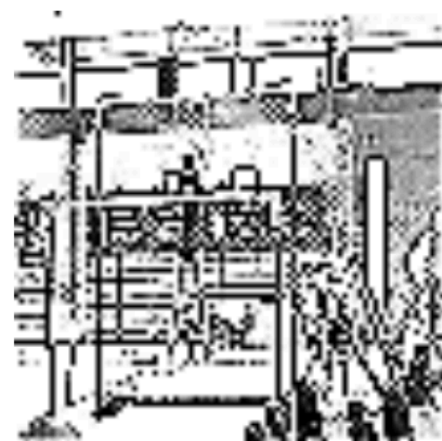
(a) plain



(b) blur

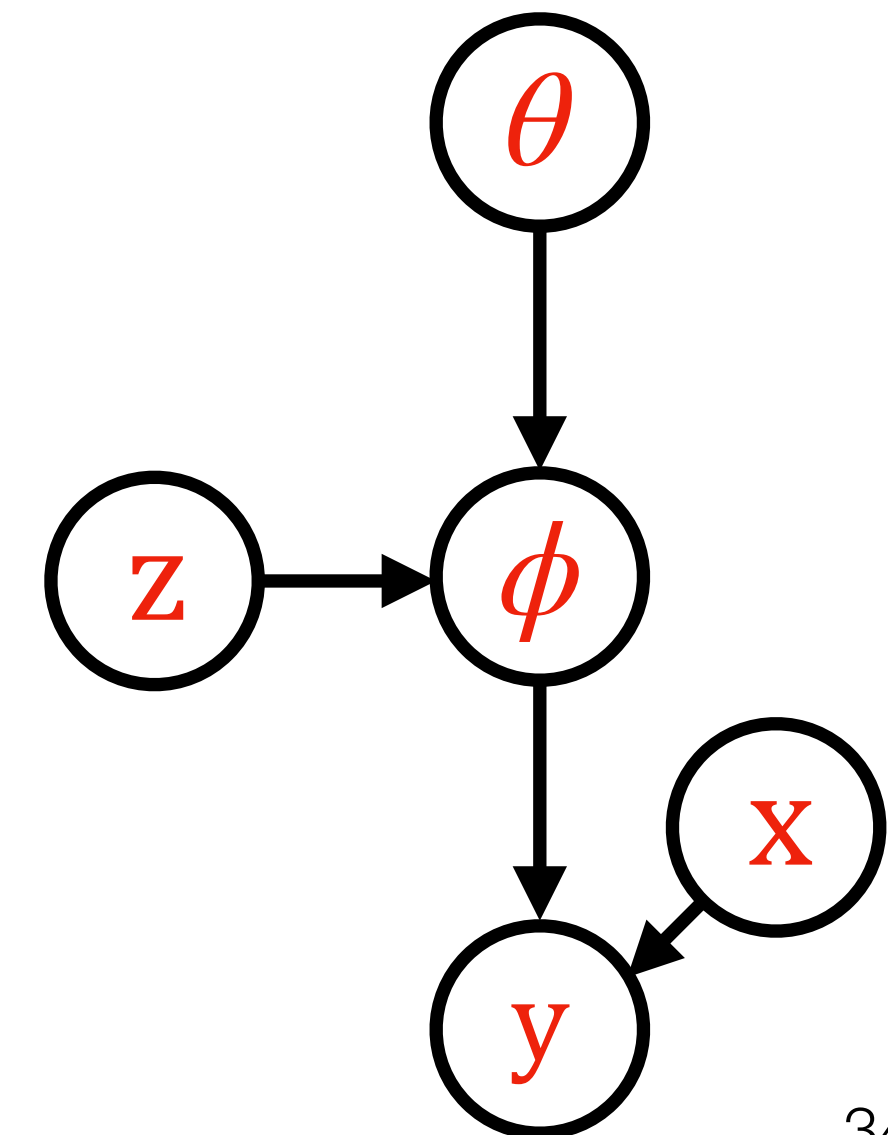


(c) night



(d) pencil

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



Meta-learning optimizers

¹ Bengio et al. 1995

² Schmidhuber 1992

³ Runarsson and Jonsson 2000

⁴ Hochreiter 2001

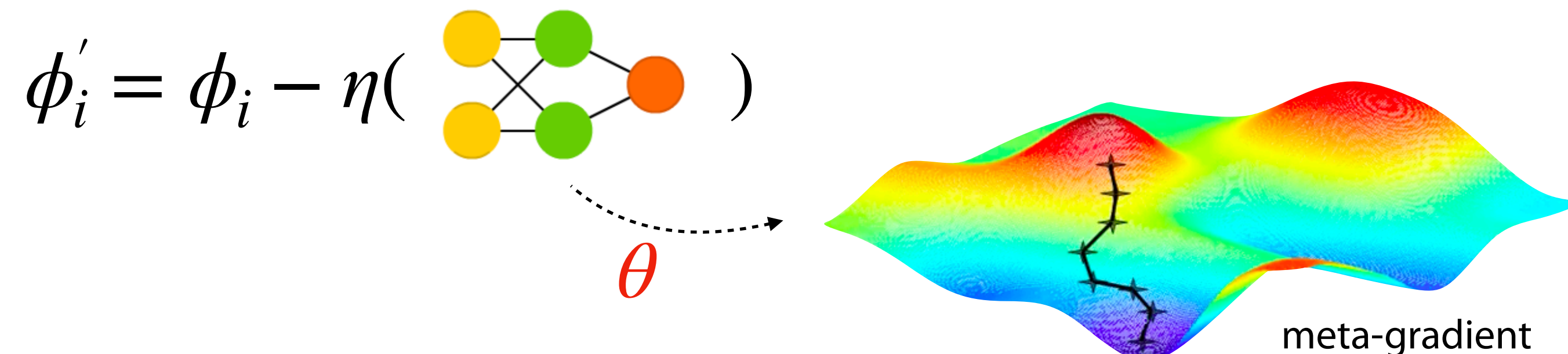
- Our brains *probably* don't do backprop, instead:

- Simple bio-inspired rules to update weights ¹

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

Diagram illustrating the update rule with annotations: η is labeled "learning rate", $y_{pre(i)}$ is labeled "presynaptic activity", and k is labeled "reinforcing signal".

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



Meta-learning optimizers

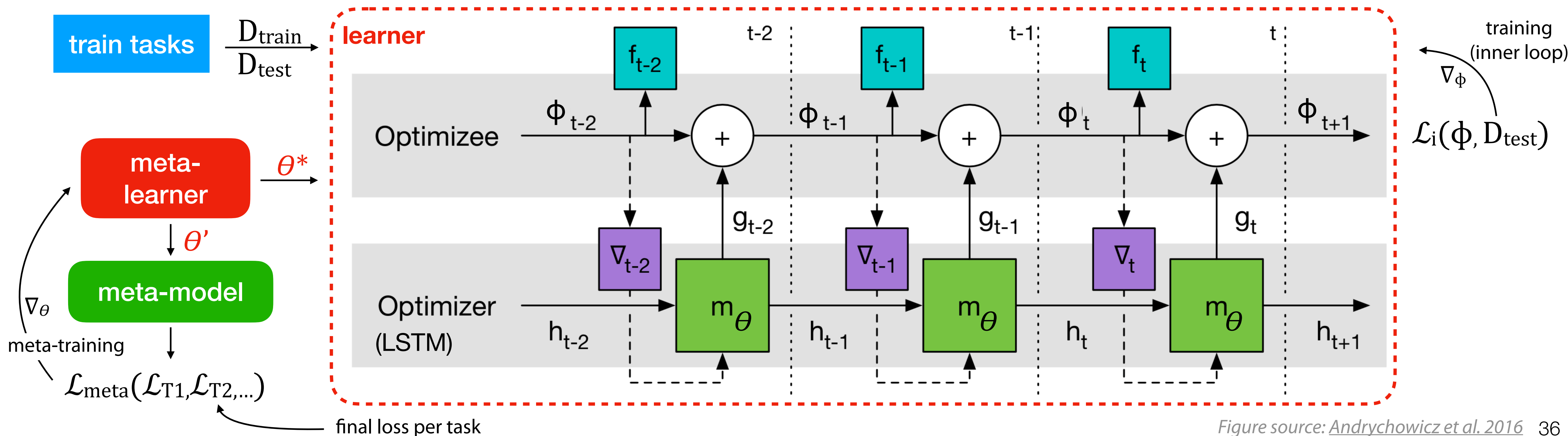
¹ Hochreiter 2001

² Andrychowicz et al. 2016

³ Ravi and Larochelle 2017

⁴ Wichrowska et al. 2017

- Represent update rule as an LSTM ^{1,2,3}, hierarchical RNN ⁴
 - Optimizee: receives weight update g_t from meta-learned optimizer
 - Optimizer: receives gradient estimate ∇_t from optimizee
- Meta-learner learns optimizer parameters with gradient descent across tasks
 - e.g. Image classification ^{2,4}, few-shot learning ³



Meta-learning optimizers

¹[Maheswaranathan et al. 2020](#)

²[Ke and Malik 2016](#)

³[Li et al. 2017](#)

⁴[Antoniou et al. 2019](#)

⁵[Park et al. 2019](#)

⁶[Flennerhag et al. 2019](#)

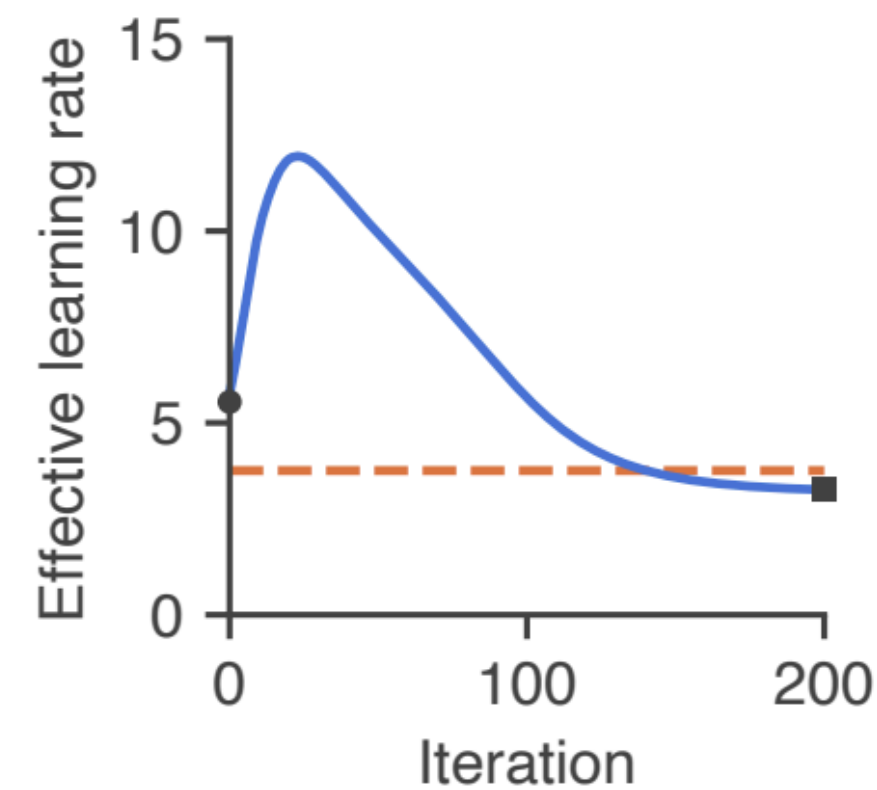
⁷[Chen et al. 2017](#)

⁸[Souza et al. 2020](#)

⁹[Kirsch et al. 2020](#)

- Meta-learned (RNN) optimizers 'rediscover' momentum, gradient clipping, learning rate schedules, learning rate adaptation, ... ¹

learning rate schedules



learning rate adaptation

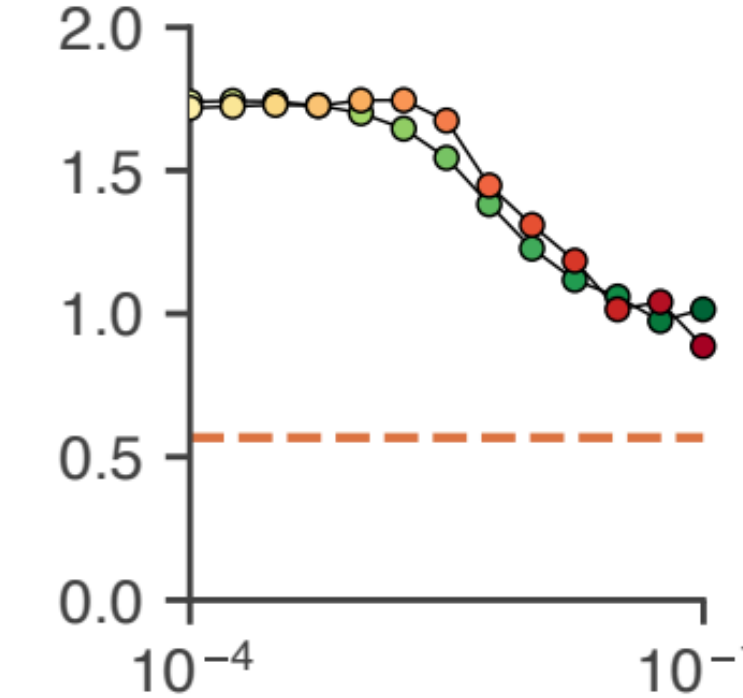
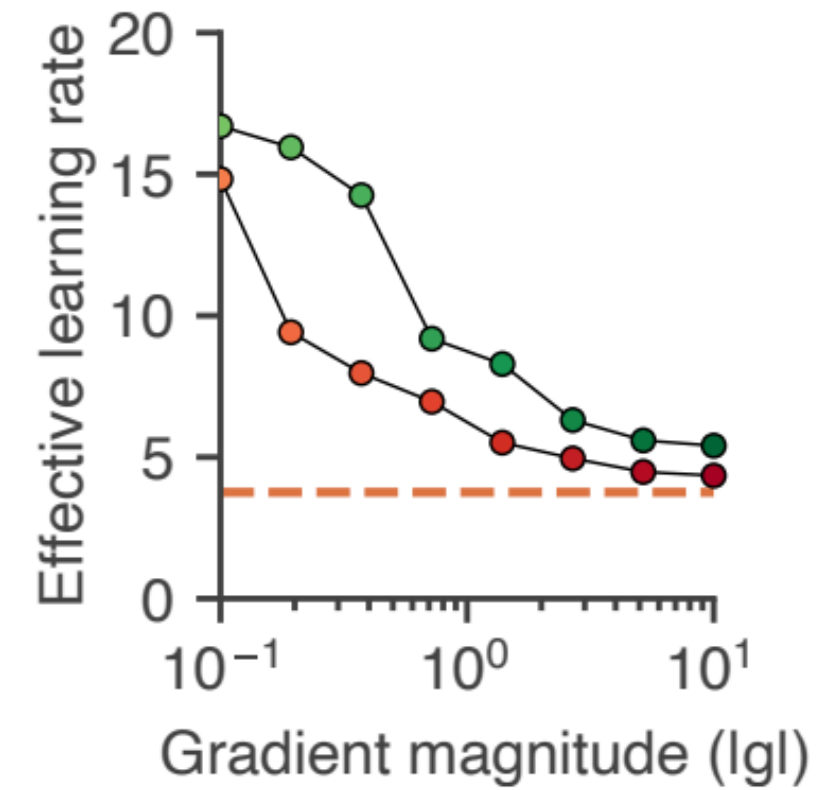
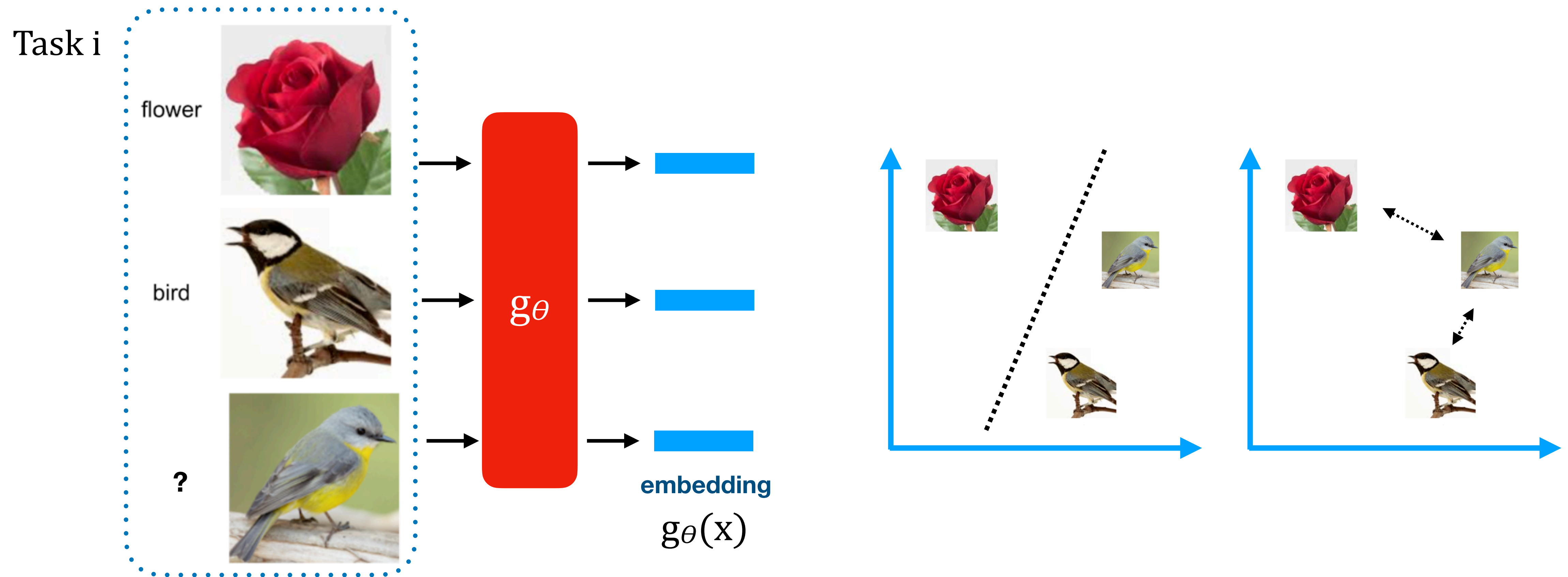


Figure source:
[Maheswaranathan et al. 2020](#)

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

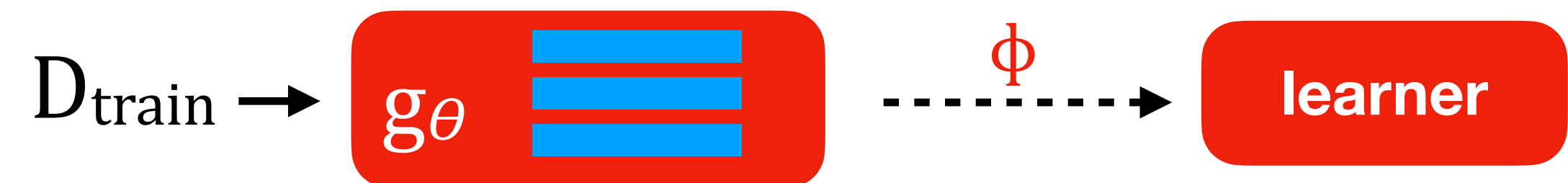
Metric learning

Learn an embedding network θ that transforms data $\{D_{\text{train}}, D_{\text{test}}\}$ 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 ϕ)



Matching networks

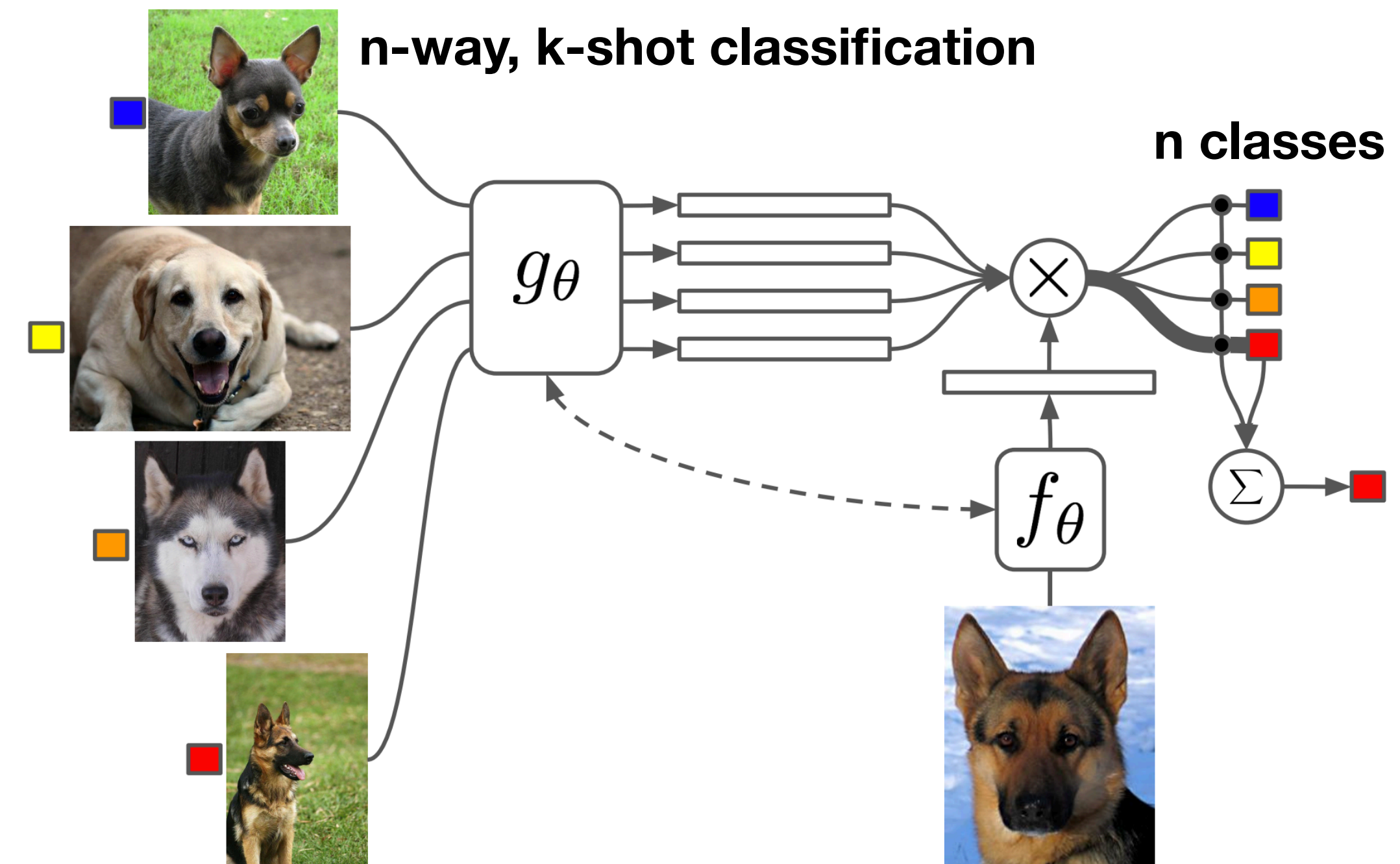
- Classifier: probability distribution based on similarity between x_{test} and $x_i \in D_{train}$

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

- 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_{θ} for D_{train} samples
 - f_{θ} 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_θ to encode each data point
- Define a prototype v_c for every class c based on the examples of that class $D_{train}^{(y)}$

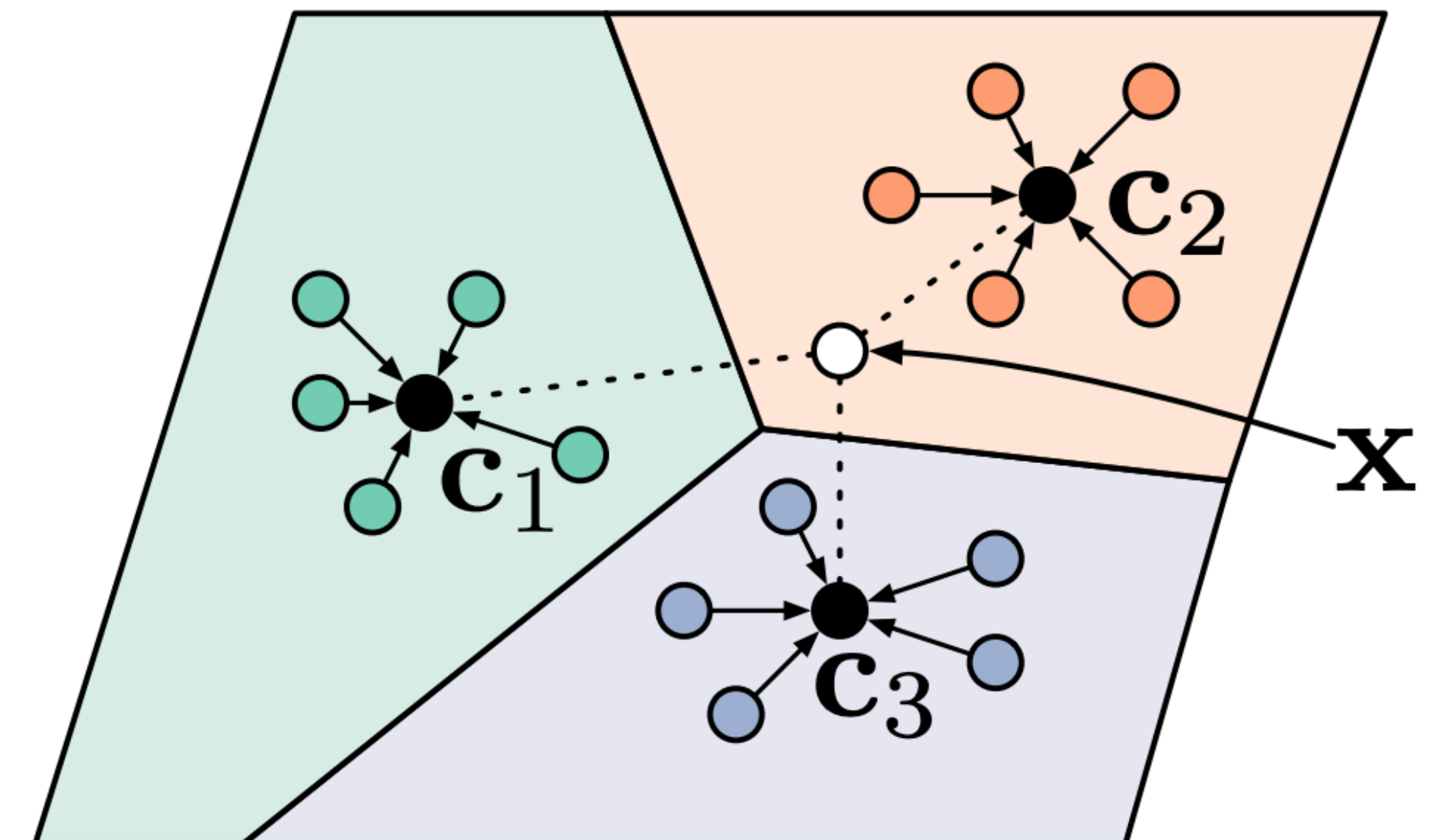
$$v_c = \frac{1}{|D_{train}^{(y)}|} \sum_{(x_i, y_i) \in D_{train}^{(y)}} f_\theta(x_i)$$

- 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:

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



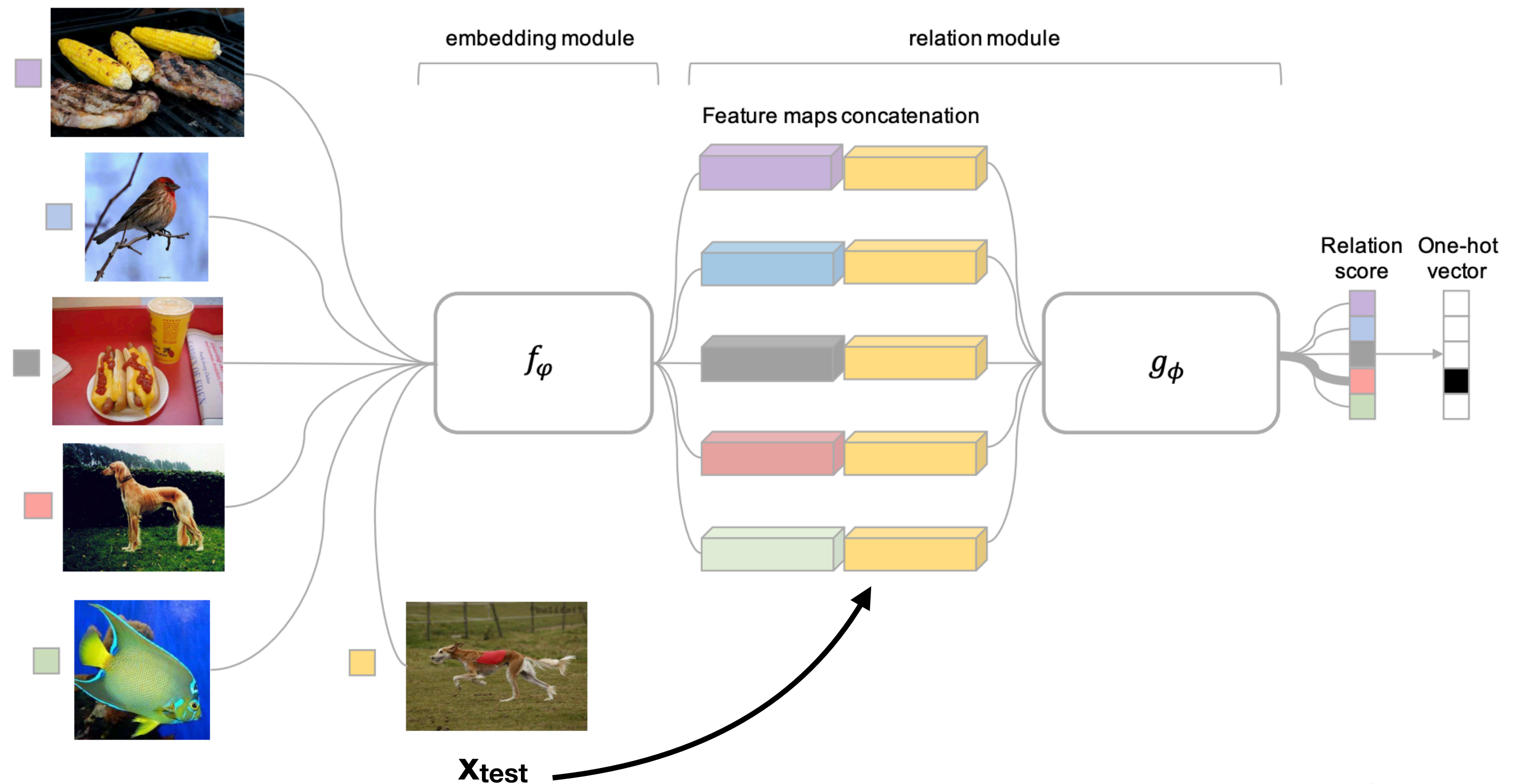
Relation networks

- Similar to matching networks, but with a trainable similarity metric
 - Learns a non-linear relationship between the data points
- Learn an embedding network f_θ and a relation network g_θ

- The relationship/similarity between a pair of inputs:

$$r_{ij} = g_\theta(\text{concat}(f_\theta(x_i), f_\theta(x_j)))$$

- Predictions based on the examples most related to x_{test}
- More expressive power, often beats other metric learners



Metric learning

Koch et al. 2015

Garcia et al. 2017

Shyam et al. 2017

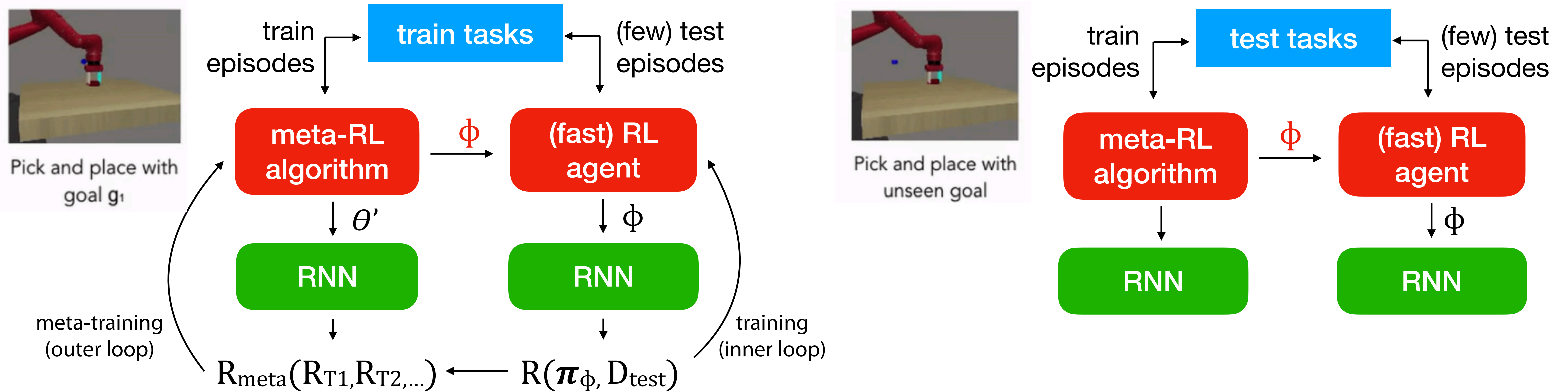
Lee et al. 2019

- Quite a few other techniques exist
 - Siamese neural networks ¹
 - Graph Neural Networks ²
 - Also applicable for semi-supervised and active learning
 - Attentive Recurrent Comparators ³
 - Compares inputs not as a whole but by parts (e.g. image patches)
 - MetaOptNet ⁴
 - 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

Wang et al. 2016

Duan et al. 2018



Previous inputs stored in hidden state

Reset the hidden state

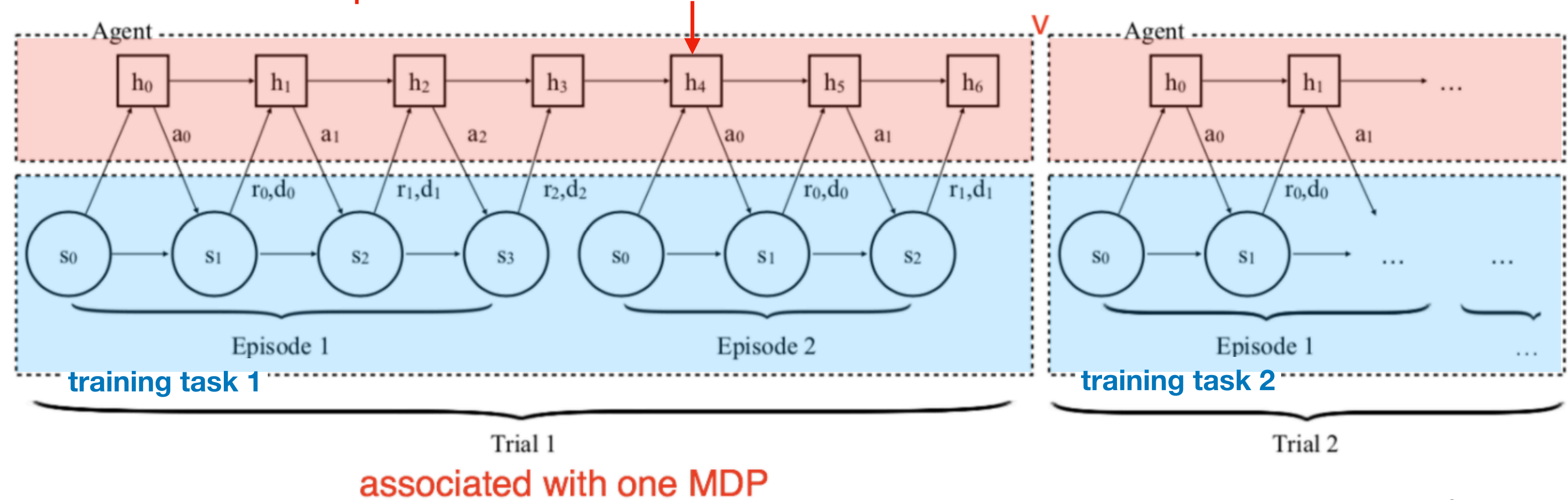


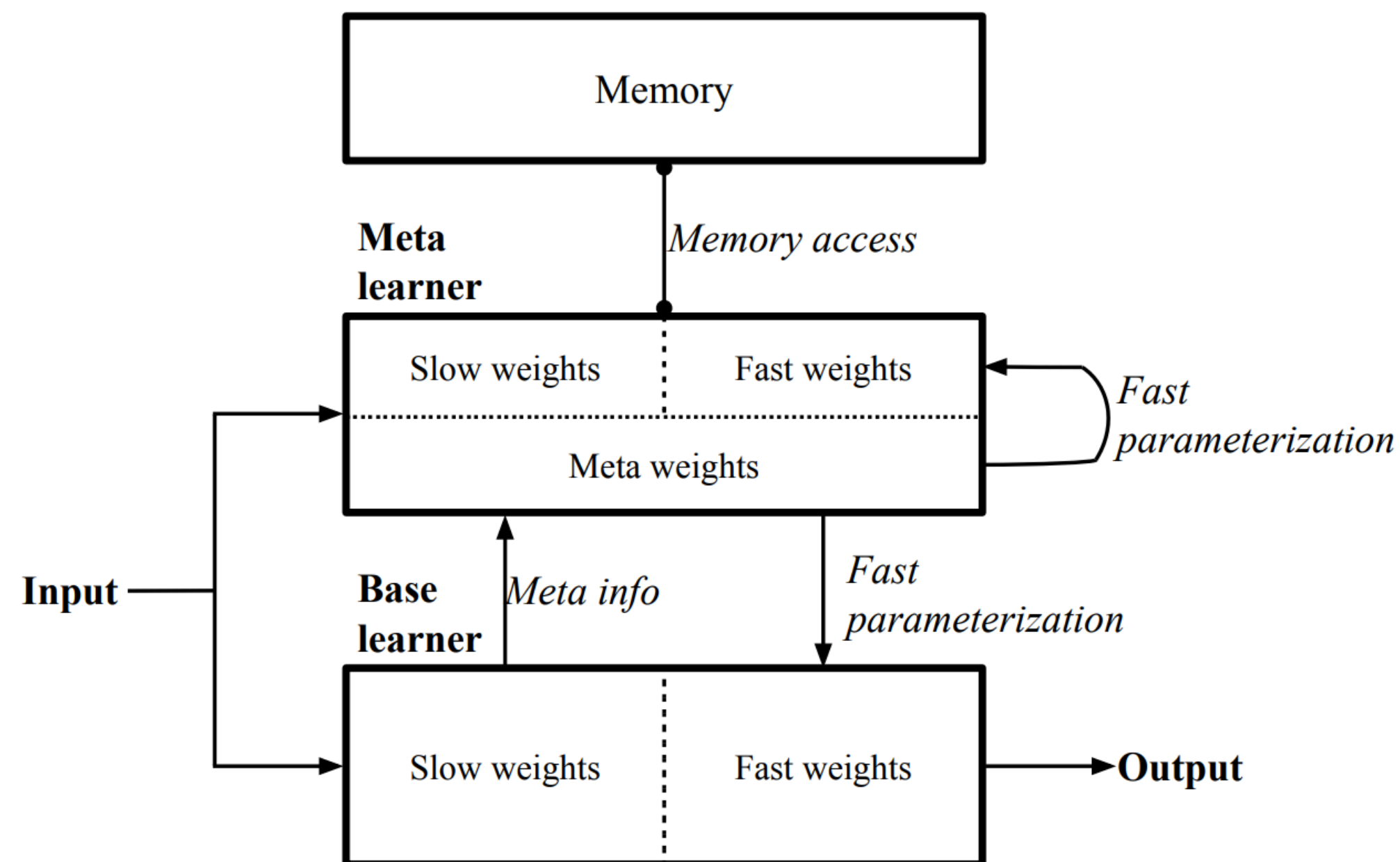
Image source: Duan et al. 2018 43

- **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

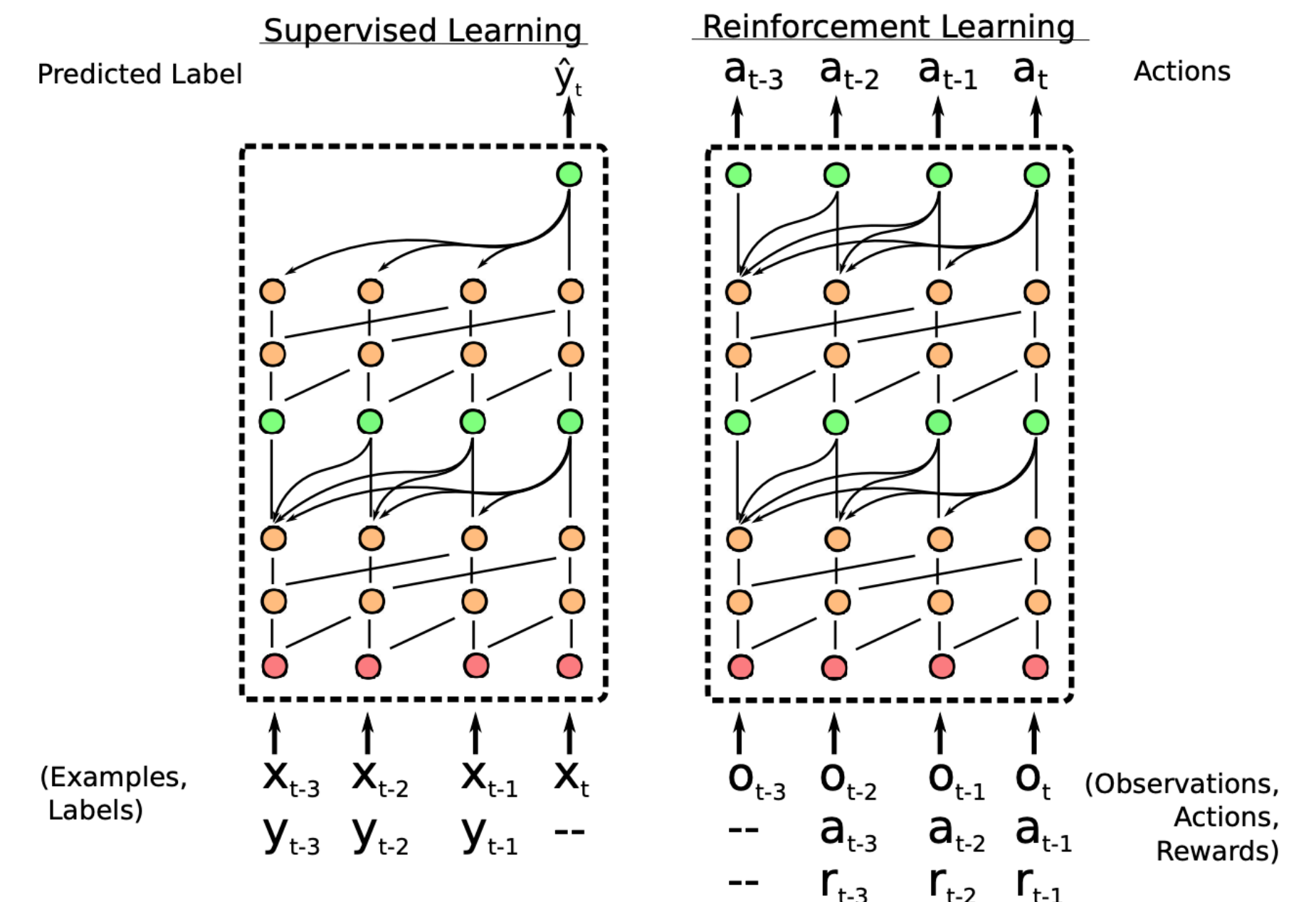
Other black-box models

Santoro et al. 2016
Munkhdalai et al. 2017
Mishra et al. 2018

- Memory-augmented NNs ¹
 - Uses neural Turing machines: short term + long term memory
- Meta Networks ²
 - 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



Meta Networks. Image source: Munkhdalai et al. 2017



SNAIL. Image source: Mishra et al. 2018

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 θ_E for agent θ_A**
 - Select agents that can solve challenges AND evolve environments so they are solvable

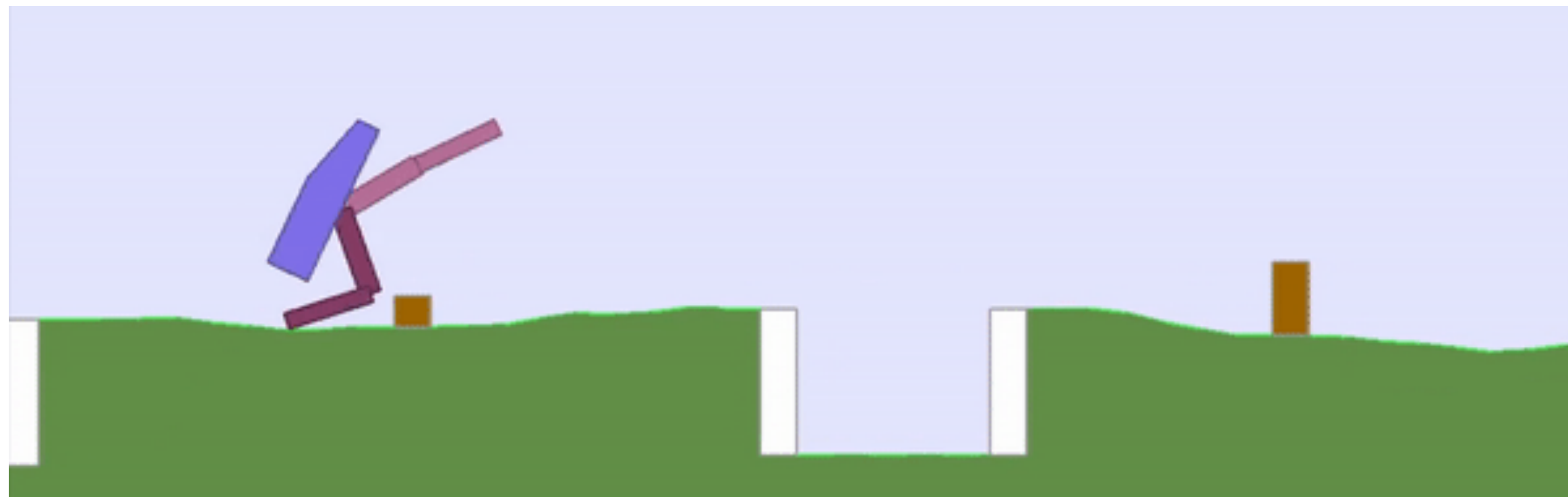
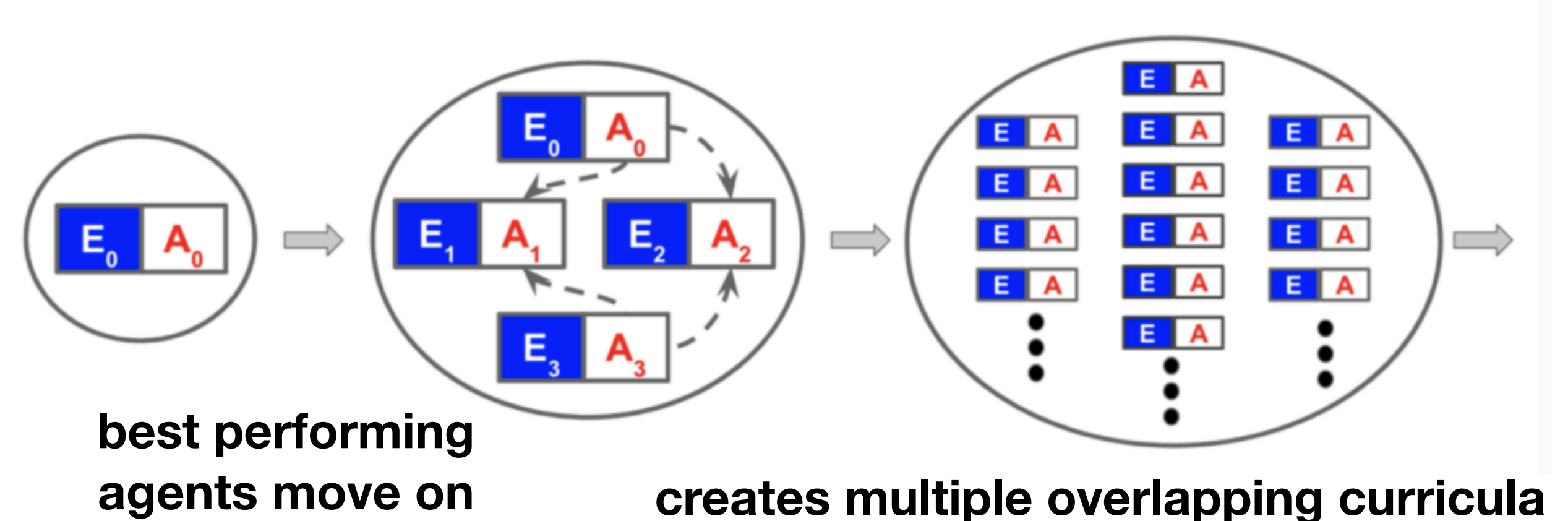


Figure source: Wang et al. 2019



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 θ_E for agent θ_A**
 - Select agents that can solve challenges AND evolve environments so they are solvable

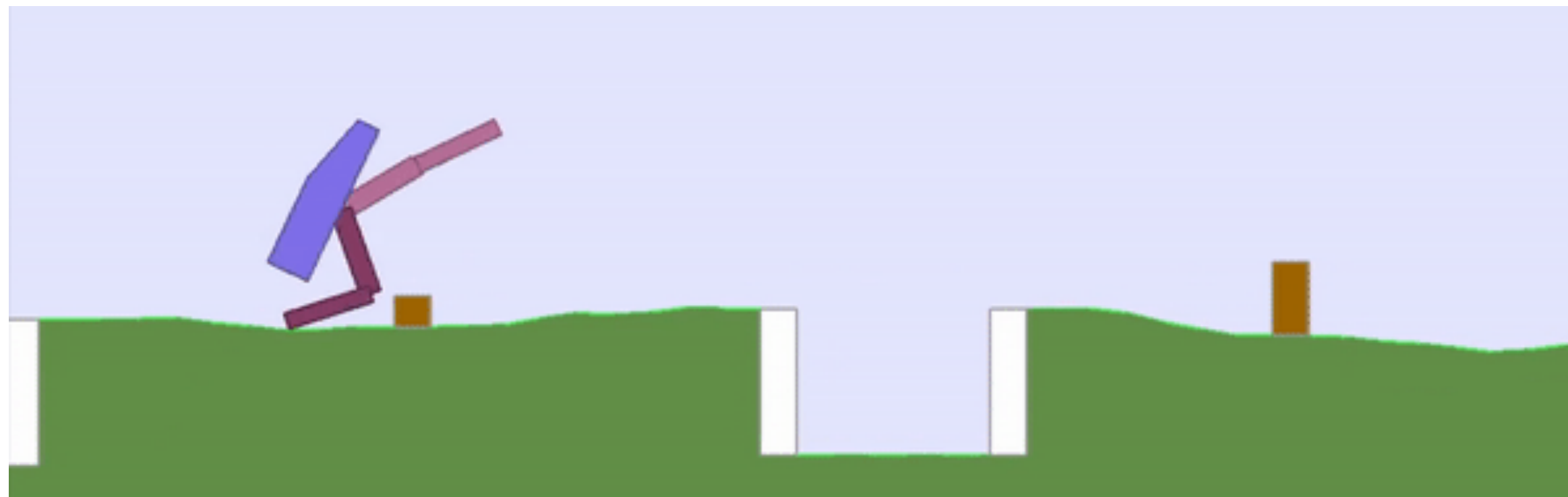
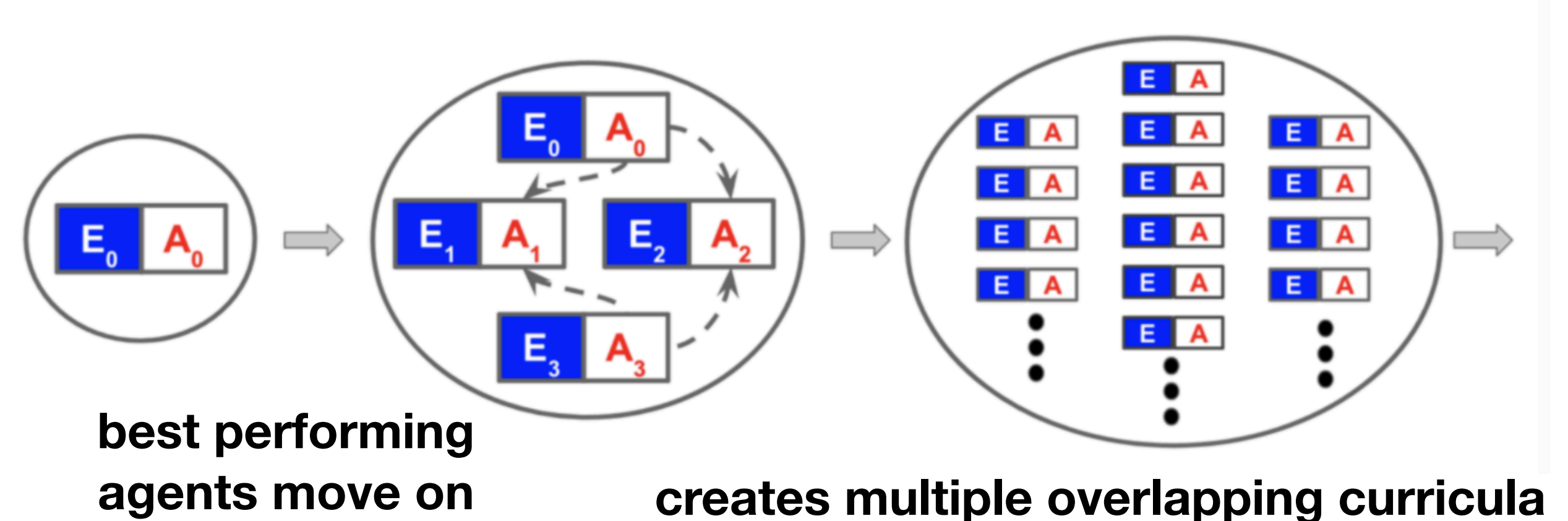


Figure source: Wang et al. 2019



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



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



Thank you

