# Learning to reinforcement learn for Neural Architecture Search *
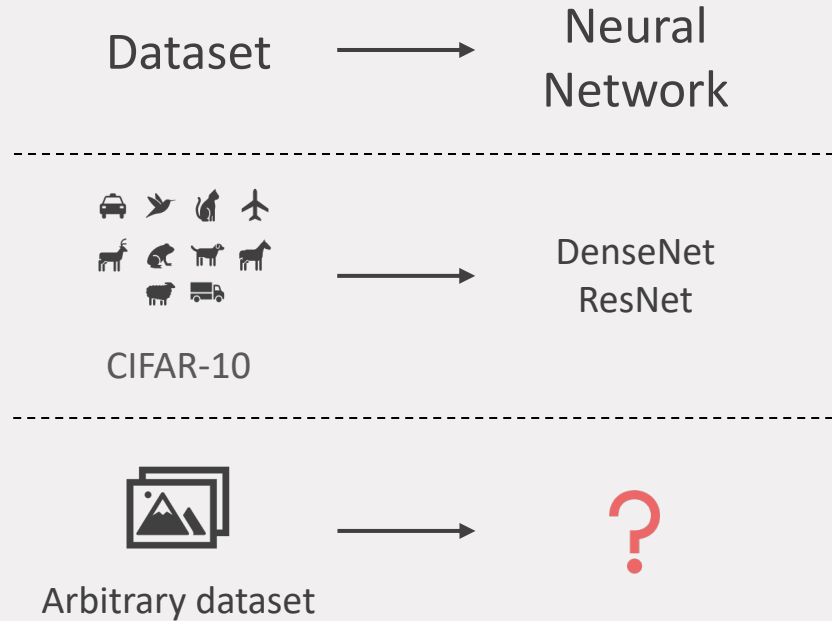
**Jorge Gómez Robles**

* M.Sc. project supervised by dr. ir. Joaquin Vanschoren

Department of Mathematics and Computer Science

# The Neural Architecture Search problem

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Neural Networks for Image Classification

Dataset $\longrightarrow$ Neural Network

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

CIFAR-10 $\longrightarrow$ DenseNet ResNet

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Arbitrary dataset $\longrightarrow$ ?

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Neural Architecture Search (NAS)

## Goal

- Automate the creation of neural architectures for **any** dataset of interest.

## Popular methods

- Bayesian optimization
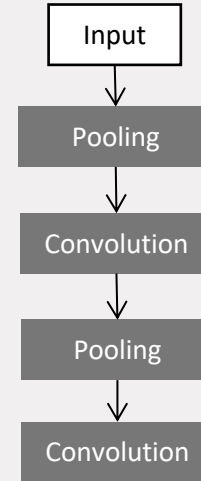- Evolutionary algorithms
- Reinforcement learning (RL)

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# The current NAS with RL framework



CIFAR-10 dataset

Reinforcement Learning

Input

Pooling

Convolution

Pooling

Convolution

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e
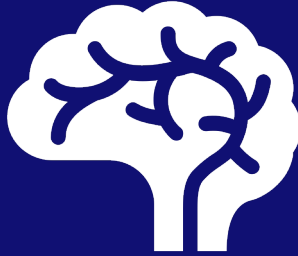
# Flaws of the current NAS framework

- The run is very expensive
  - 28 days with 500 GPUs (B. Zoph and Q. V Le, "Neural Architecture Search with Reinforcement Learning", 2017)
  - 3 days with 32 GPUs (Z. Zhong *et al.*, "BlockQNN: Efficient Block-wise Neural Network Architecture Generation", 2018)
- For every new environment (i.e., dataset), a RL run must be performed from scratch.

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# The research question

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles
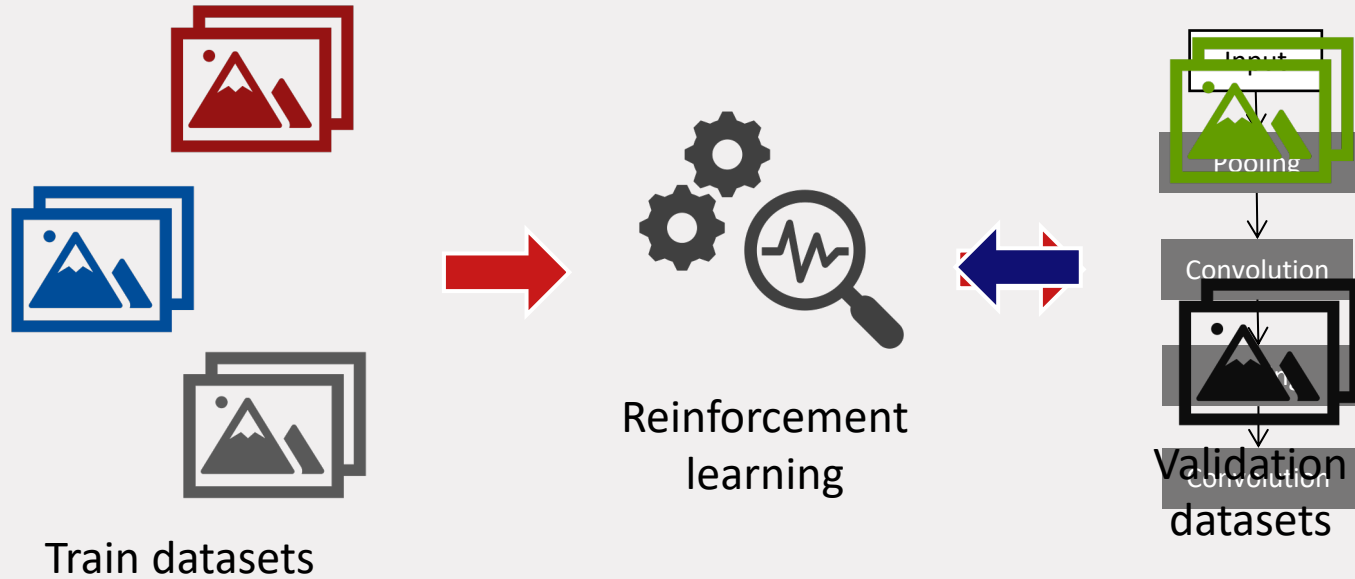
TU/e

# Research question

- Is it possible to learn an **adaptive strategy (policy)** to design CNNs for image classification so that we can transfer it to avoid **training from scratch**?
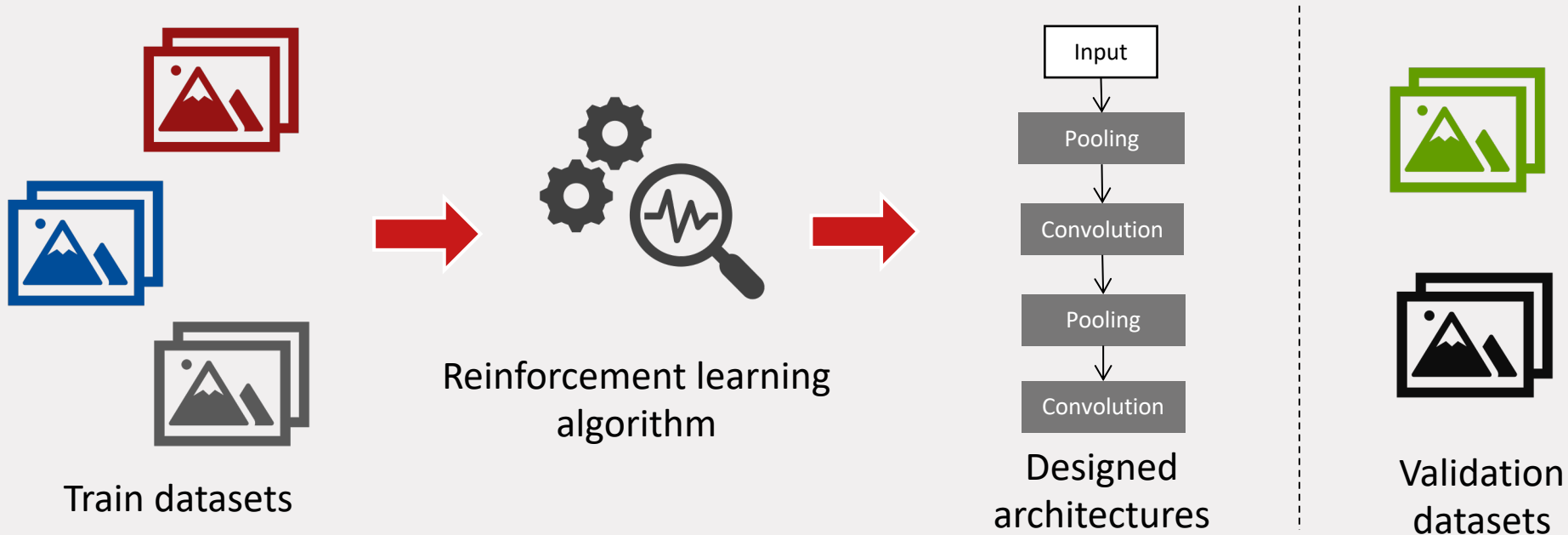
Consequences:

1. Transfer the policy and use it as an initialization to train on new environments
2. Transfer the policy to new environments without training

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Proposed solution



"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

Reinforcement learning

Train datasets

Validation datasets

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Analyzing the components thoroughly



Train datasets

Reinforcement learning algorithm

Input

Pooling

Convolution

Pooling

Convolution

Designed architectures

Validation datasets

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Analyzing the components thoroughly

Use a **meta-reinforcement learning algorithm:**

- J. X. Wang et al., "*Learning to reinforcement learn*", 2016.
- Build on top of a the A2C algorithm
- An LSTM-based meta-learner that learns the relation between all the agent-environment interactions within an episode.

Reinforcement learning

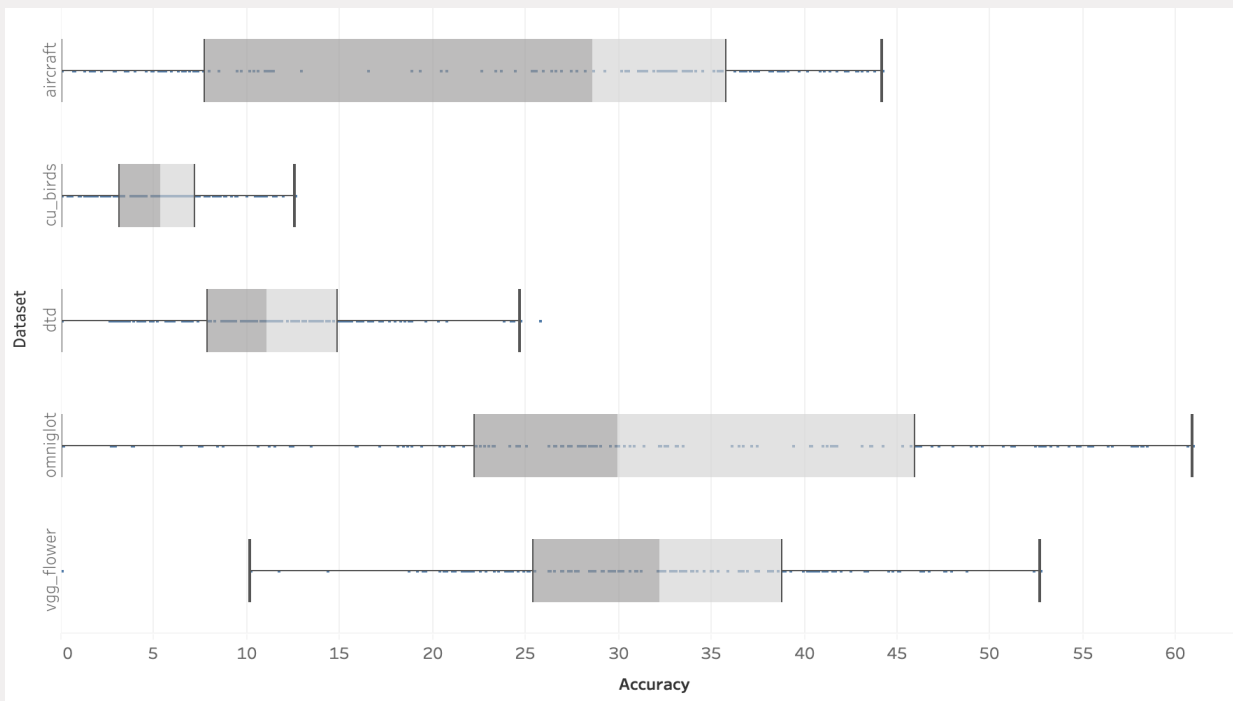"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Analyzing the components thoroughly

Five datasets from the the **meta-dataset**.

- E. Triantafillou et al., *"Meta-Dataset : A Dataset of Datasets for Learning to Learn from Few Examples"*, 2018.

Datasets

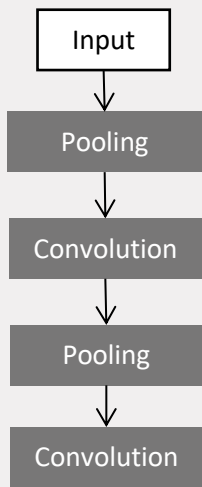| Dataset | ID | N Classes | N Obs | Usage |
|---|---|---|---|---|
| FGVC-Aircraft | aircraft | 100 | 10,000 | Validation |
| CUB-200-2011 | cu_birds | 200 | 11,788 | Validation |
| Describable Textures | dtd | 47 | 5,640 | Training |
| VGG Flower | vgg_flower | 102 | 8,189 | Training |
| Omniglot | omniglot | 1623 | 32,460 | Training |

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

# Analyzing the components thoroughly



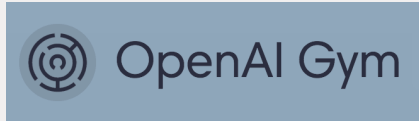Designed architectures

Based on the methodology of:

- Z. Zhong et al., *"Practical Block-wise Neural Network Architecture Generation"*, 2017.

Highlights:

- 6 types of layers in the network
- Chain-structured networks only
- Maximum number of layers: 10
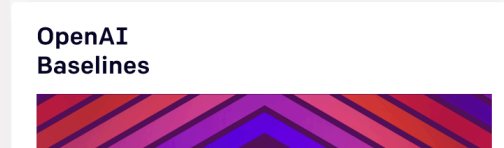- Early-stop training to avoid long runs (12 epochs)

TU/e

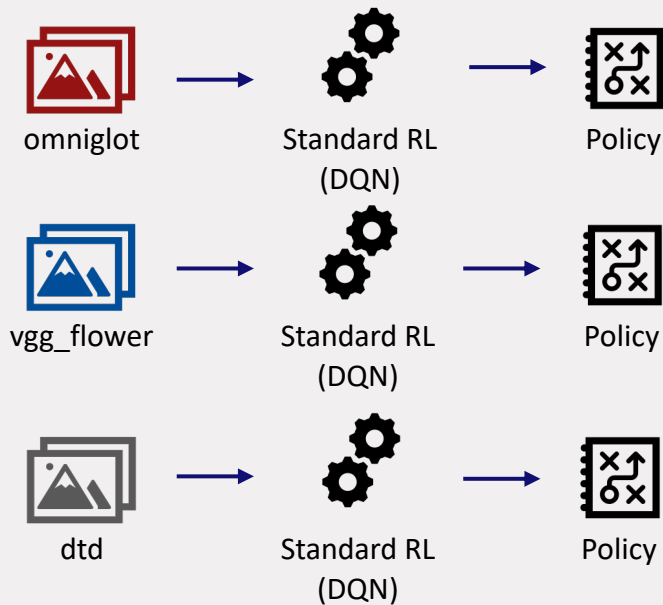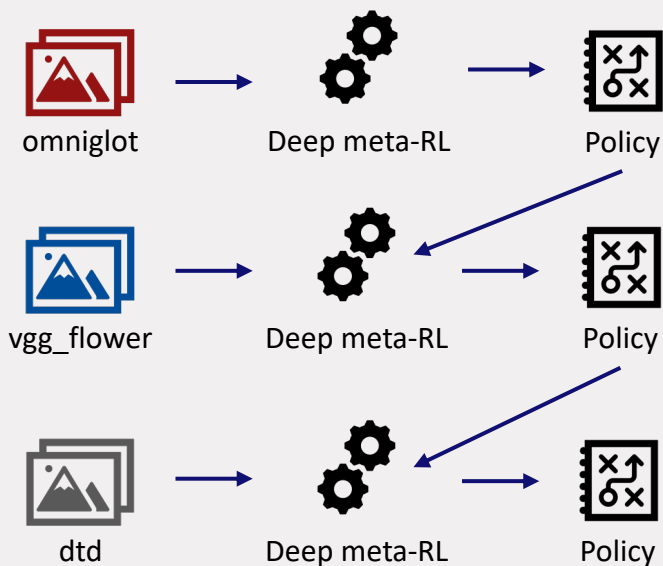| Type of layer | Hyper-parameter |
| --- | --- |
| Convolution | Kernel size = {1, 3, 5} |
| MaxPooling | Pool size = {2, 3} |
| AvgPooling | Pool size = {2, 3} |
| Terminal | - |

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Software



NasGym
github.com/gomerudo/nas-env



Deep meta-RL
github.com/gomerudo/openai-baselines

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Experiments and results

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

# Experiment 1: training with deep meta-RL



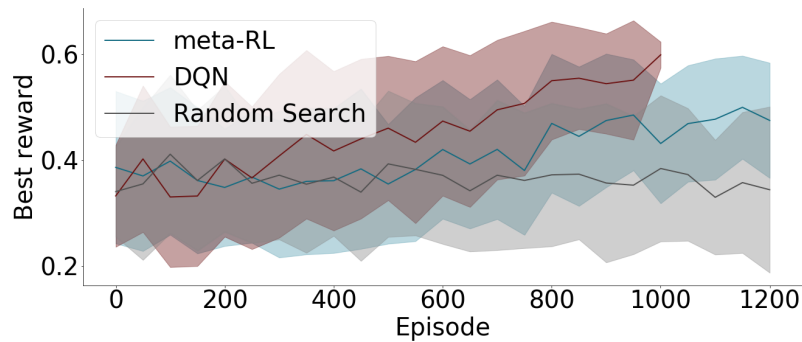"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles
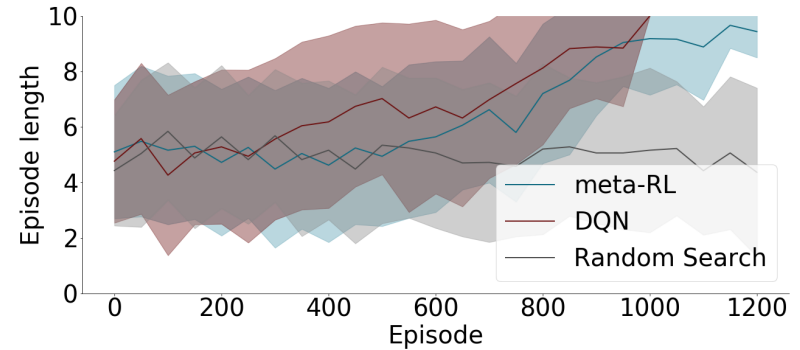
Goal:

- Observe if deep meta-reinforcement learning gives an advantage over standard RL
    - *Best reward* -> Quality (the best designed network)
    - *Episode length* -> Complexity (episode length equal to the depth of the network)
- Study if the agent can adapt
    - *Policy entropy* -> exploration
    - *Proportion of actions* -> strategy per environment

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

**TU/e**

# Results of Experiment 1: omniglot



Avg. 'Best reward' every 50 episodes for 'omniglot'
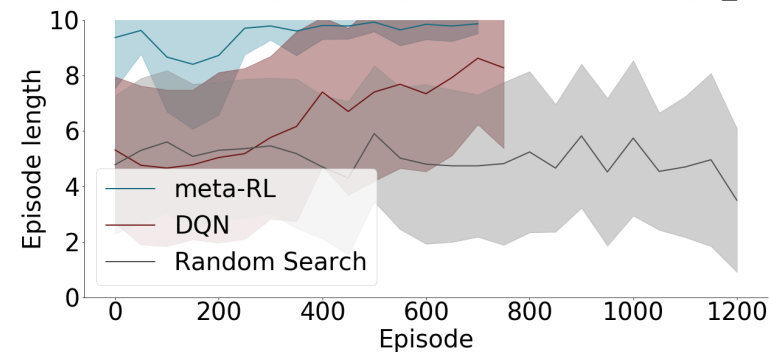
Avg. 'Episode length' every 50 episodes for 'omniglot'

\* Confidence range is 1 STD

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Results of Experiment 1: vgg_flower
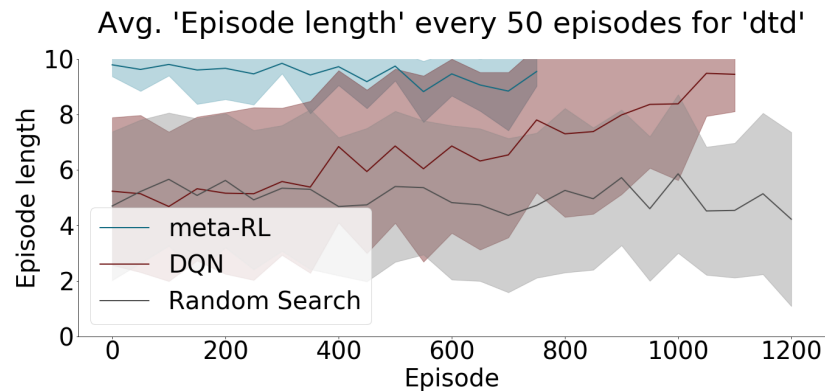


Avg. 'Best reward' every 50 episodes for 'vgg_flower'



Avg. 'Episode length' every 50 episodes for 'vgg_flower'

\* Confidence range is 1 STD

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Results of Experiment 1: dtd



Avg. 'Best reward' every 50 episodes for 'dtd'



Avg. 'Episode length' every 50 episodes for 'dtd'

* Confidence range is 1 STD

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Results of Experiment 1: exploration per environment



Evolution of the 'Policy entropy' for deep meta-RL

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

# Results of Experiment 1: the strategy



Proportion of actions performed on the training datasets

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

# Experiment 2: evaluating the policy (Part A)



"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

Goal:

- Study the strategies deployed by the agent

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Results of Experiment 2: aircraft



Avg. 'Best reward' every 50 episodes for 'aircraft'



Avg. 'Episode length' every 50 episodes for 'aircraft'

* Confidence range is 1 STD

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Results of Experiment 2: cu_birds



Avg. 'Best reward' every 50 episodes for 'cu_birds'

Avg. 'Episode length' every 50 episodes for 'cu_birds'

* Confidence range is 1 STD

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Results of Experiment 2: the strategy



Proportion of actions executed on the evaluation datasets

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Experiment 2: evaluating the policy (Part B)



aircraft

Best 2 Neural Architectures

cu_birds

Best 2 Neural Architectures

Intensive Training (100 epochs)

Shorthened VGG19

Intensive Training (100 epochs)

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

| Dataset | Deep meta-RL (1st) | Deep meta-RL (2nd) | Shortened VGG19 |
|---|---|---|---|
| aircraft | $49.18 \pm 1.20$ | **$50.11 \pm 1.02$** | $30.85 \pm 10.82$ |
| cu_birds | $23.97 \pm 1.28$ | **$24.24 \pm 0.90$** | $6.66 \pm 1.98$ |

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Conclusions

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Conclusions

- Deep meta-RL shows a good behavior when the policy is transferred during training
  - Shows indications of adaptation between environments
  - Outperforms standard RL
  - Shows consistency (small variance)
- During the evaluation the behavior can potentially be improved
  - The strategies deployed on different environments are not *ad-hoc*
  - The agent can design networks that outperform manually engineered solutions

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles
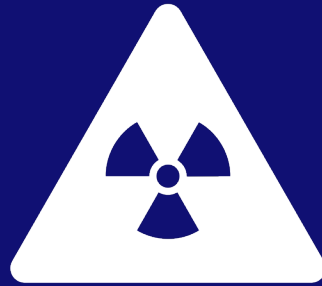
TU/e

# What can we do better?

- Hyper-parameter tuning of deep meta-RL
  - Encourage exploration
  - Make learning faster
- Increasing the number of time-steps for the agent-environment interaction during training
  - Might help to improve the results during evaluation

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Thanks!

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Back-up slides

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Results of Experiment 1: running times

| Environment (dataset) | Deep meta-RL | DQN |
|---|---|---|
| omniglog | 11 days 9h | 6 days 14h |
| vgg_flower | 7 days 23h | 5 days 15h |
| dtd | 6 days 17h | 6 days 4h |
| **Total** | **26 days 1h** | **18 days 9h** |

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Results of Experiment 2: running times

| Environment (dataset) | Time |
|---|---|
| aircraft | 2 days 6h |
| cu_birds | 2 days 22h |

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

# The deep meta-RL interaction



"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles
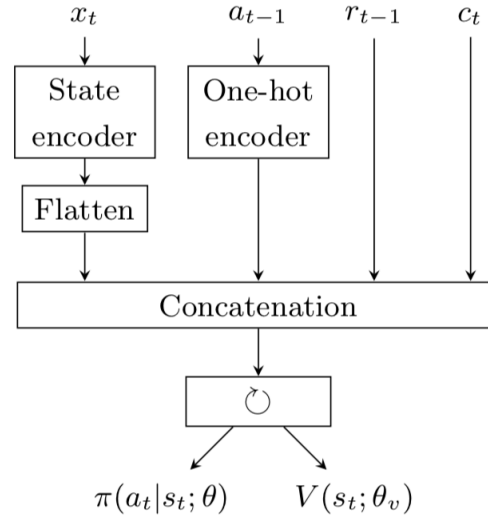
# The A2C architecture



Figure 7: Illustration of the *meta*-A2C architecture. In our implementation, the "State encoder" follows the procedure explained in Section 4.2.1, and the recurrent layer is an LSTM with 128 units.

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

# The policy optimization function

Formally, let $t$ be the current time step, $s_t = x_t \cdot a_{t-1} \cdot r_{t-1} \cdot c_t$ a concatenation of inputs, $\pi(a_t|s_t; \theta)$ the policy, $V(s_t; \theta_v)$ the value function, $H$ the entropy, $j \in \mathbb{N}$ the horizon, $\gamma \in (0, 1]$ the discount factor, $\eta$ the regularization coefficient, and $R_t = \sum_{i=0}^{j-1} \gamma^i r_{t+i}$ the total accumulated return from time step $t$. The gradient of the objective function is:

$$\nabla_\theta \log \pi(a_t|s_t; \theta) \underbrace{(R_t - V(s_t; \theta_v))}_{\text{Advantage estimate}} + \underbrace{\eta \nabla_\theta H(\pi(s_t; \theta))}_{\text{Entropy regularization}} \tag{1}$$

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

**TU/e**

# Prediction module

- Dense layer 1024 units
- ReLU
- Dropout 0.4
- Dense layer n_classes

# Training

- Learning rate: 0.001
  - Reduced by a factor of 0.2 every 5 epochs
- Beta1 = 0.9
- Beta2 = 0.999
- Epsilon= 10e-8

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Time-steps for experiment 1

| Environment | Deep meta-RL | DQN |
|---|---|---|
| omniglot | 8000 | 6500 |
| vgg_flower | 7000 | 5500 |
| dtd | 7000 | 7000 |

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e

# Hyper-parameters for experiment 1

| Environment | A2C | DQN | Training of networks |
|---|---|---|---|
| d = 10 (max layers)<br>$\tau = 10$ (Maximum length of episode) | j = 5 (number of steps before updating the network)<br>$\gamma = 0.9$ (discount factor)<br>$\eta = 0.01$ (Entropy regularization)<br>$\alpha = 0.001$ (Learning rate) | Buffer size = t_max / 2<br>Target's model batch size = 20<br>$\varepsilon$: Linear decay from 1 to 0.1<br>$\alpha = 0.0005$ (Learning rate) | Batch size = 128<br>Num. epochs = 12 |

"Learning to reinforcement learn for Neural Architecture Search" by J. Gómez Robles

TU/e