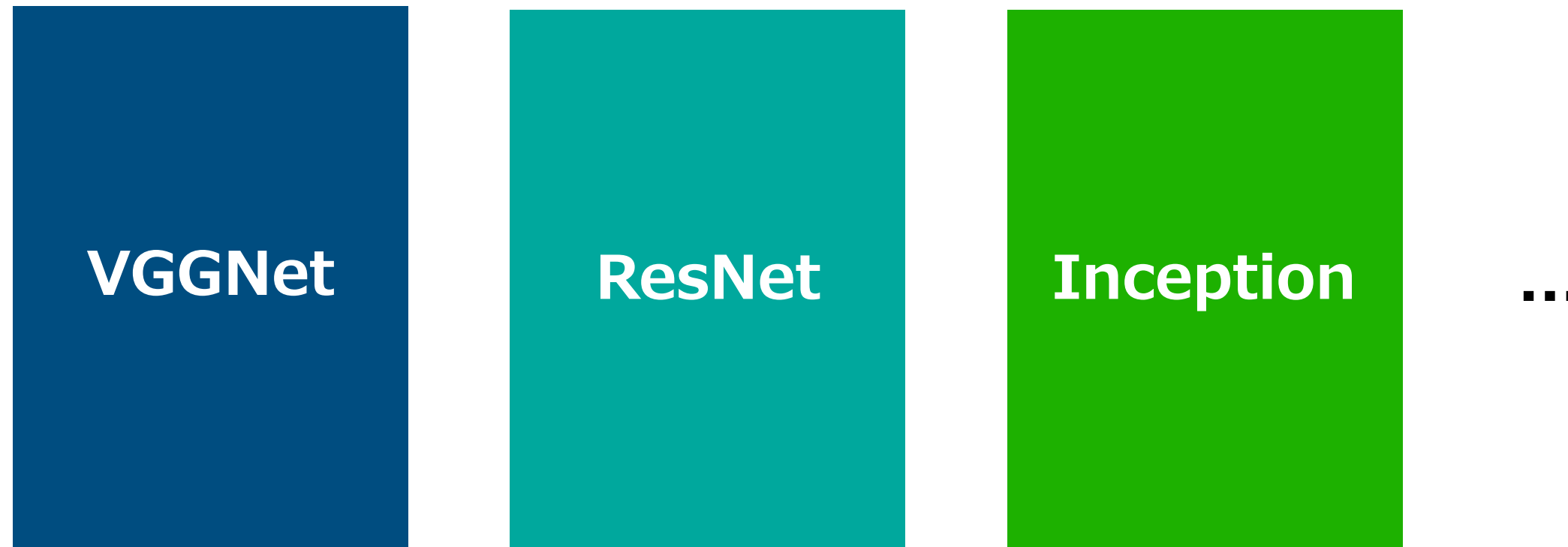


Adaptive Stochastic Natural Gradient Method for One-Shot Neural Architecture Search

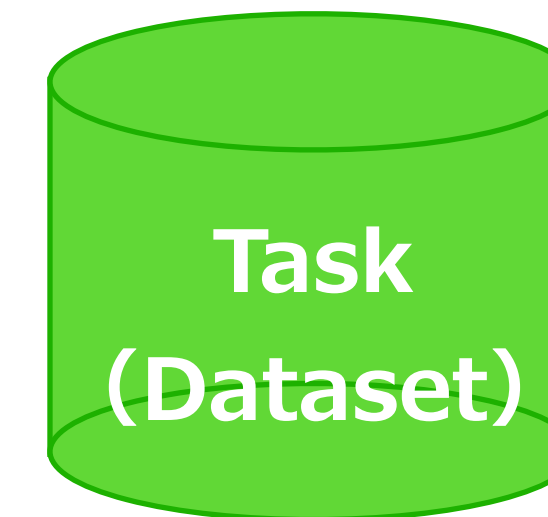
- Youhei Akimoto (University of Tsukuba / RIKEN AIP)
- Shinichi Shirakawa (Yokohama National University)
- Nozomu Yoshinari (Yokohama National University)
- Kento Uchida (Yokohama National University)
- Shota Saito (Yokohama National University)
- Kouhei Nishida (Shinshu University)

Neural Architecture

Neural Network Architectures



often pre-trained
on some datasets



Sometimes...

- a known architecture works well on our tasks.
Happy!



Other times...

- Find a good one
- Design a brand-new architecture and train it

Trial and Error!



One-Shot Neural Architecture Search

Joint Optimization of Architecture c and Weights w

NAS as hyper-parameter search

$$\max_c f(w^*(c), c)$$

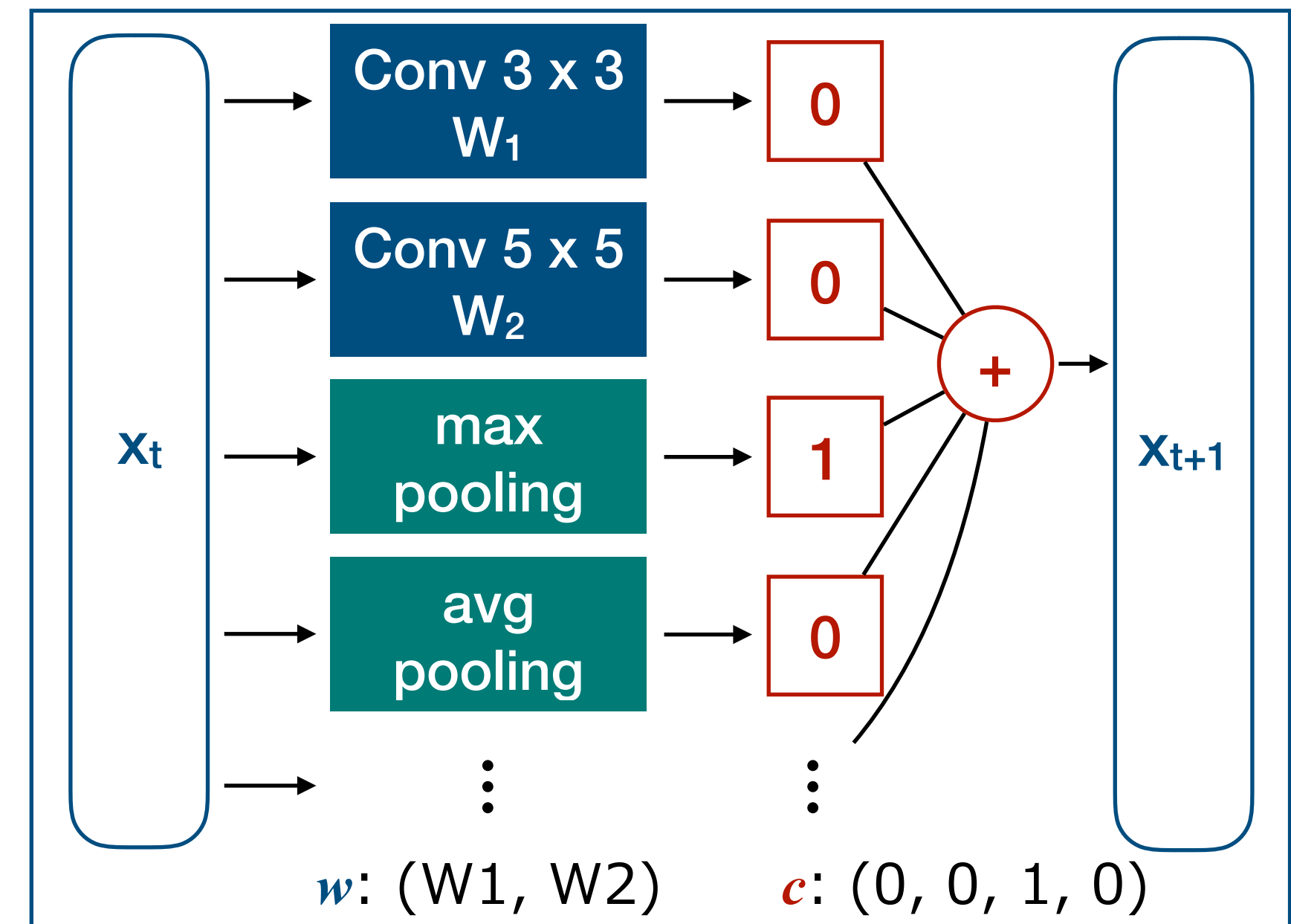
$$\text{subject to } w^*(c) = \operatorname{argmax}_w f(w, c)$$

1 c evaluation
= 1 training

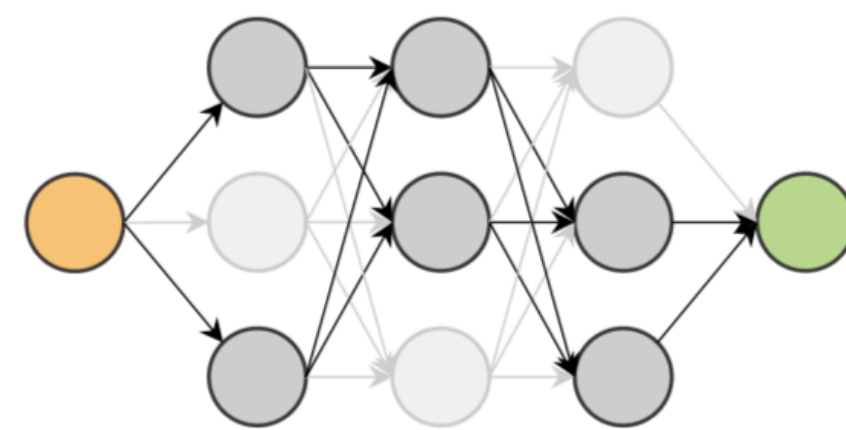
One-shot NAS

$$\max_{w, c} f(w, c)$$

optimization of x and
 c within 1 training



NAS



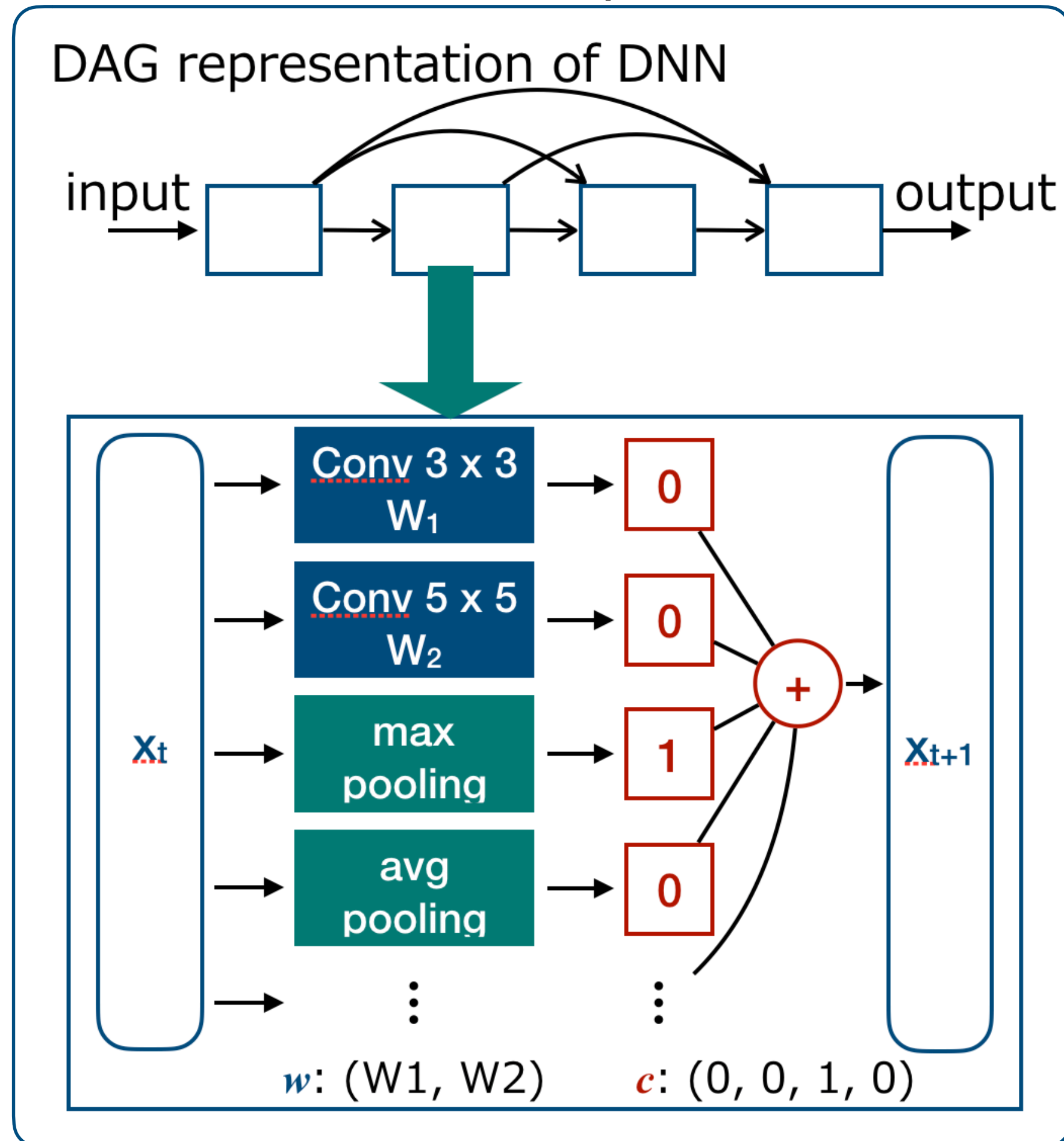
good architecture



Difficulties for Practitioners

How to choose / tune the search strategy?

Search Space



Search Strategy

Gradient-Based Method

$$w \leftarrow w + \epsilon_w \nabla_w f(w, c(\theta))$$

$$\theta \leftarrow \theta + \epsilon_\theta \nabla_\theta f(w, c(\theta))$$

hyper-parameter: **step-size**

Other Choices

- Evolutionary Computation Based
- Reinforcement Learning Based

- how to treat **integer variables** such as #filters?
- how to tune the **hyper-parameters** in such situations?

Contributions

Novel Search Strategy for One-shot NAS

1. arbitrary search space (categorical + ordinal)
2. robust against its inputs (hyper-param. and search space)

Our approach 1. Stochastic Relaxation

$$\max_{\mathbf{w}, \mathbf{c}} f(\mathbf{w}, \mathbf{c}) \Rightarrow \max_{\mathbf{w}, \boldsymbol{\theta}} \boxed{J(\mathbf{w}, \boldsymbol{\theta})} := \int f(\mathbf{w}, \mathbf{c}) \overset{\text{exponential family}}{\boxed{p(\mathbf{c} | \boldsymbol{\theta})}} d\mathbf{c}$$

differentiable w.r.t. \mathbf{w} and $\boldsymbol{\theta}$

2. Stochastic Natural Gradient + Adaptive Step-Size

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \epsilon_{\mathbf{w}}^t \nabla_{\mathbf{w}} \widehat{J}(\mathbf{w}^t, \boldsymbol{\theta}^t)$$

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t + \epsilon_{\boldsymbol{\theta}}^t \boxed{\mathbf{F}(\boldsymbol{\theta}^t)^{-1} \nabla_{\boldsymbol{\theta}} \widehat{J}(\mathbf{w}^{t+1}, \boldsymbol{\theta}^t)} \text{ Natural Gradient}$$

➔ Under appropriate step-size

$$J(\mathbf{w}^t, \boldsymbol{\theta}^t) < J(\mathbf{w}^{t+1}, \boldsymbol{\theta}^t) < J(\mathbf{w}^{t+1}, \boldsymbol{\theta}^{t+1}) \text{ Monotone Improvement}$$

Results and Details

- Faster & Competitive Accuracy to other one-shot NAS

Table 1: Comparison of different architecture search methods on CIFAR-10. The search cost indicates GPU days for architecture search excluding the retraining cost.

Method	Search Cost (GPU days)	Params (M)	Test Error (%)
NASNet-A (Zoph et al., 2018)	1800	3.3	2.65
NAONet (Luo et al., 2018)	200	128	2.11
ProxylessNAS-G (Cai et al., 2019)	4	5.7	2.08
SMASHv2 (Brock et al., 2018)	1.5	16.0	4.03
DARTS second order (Liu et al., 2019)	4	3.3	2.76 (± 0.09)
DARTS first order (Liu et al., 2019)	1.5	3.3	3.00 (± 0.14)
SNAS (Xie et al., 2019)	1.5	2.8	2.85 (± 0.02)
ENAS (Pham et al., 2018)	0.45	4.6	2.89
ASNG-NAS	0.11	3.9	2.83 (± 0.14)

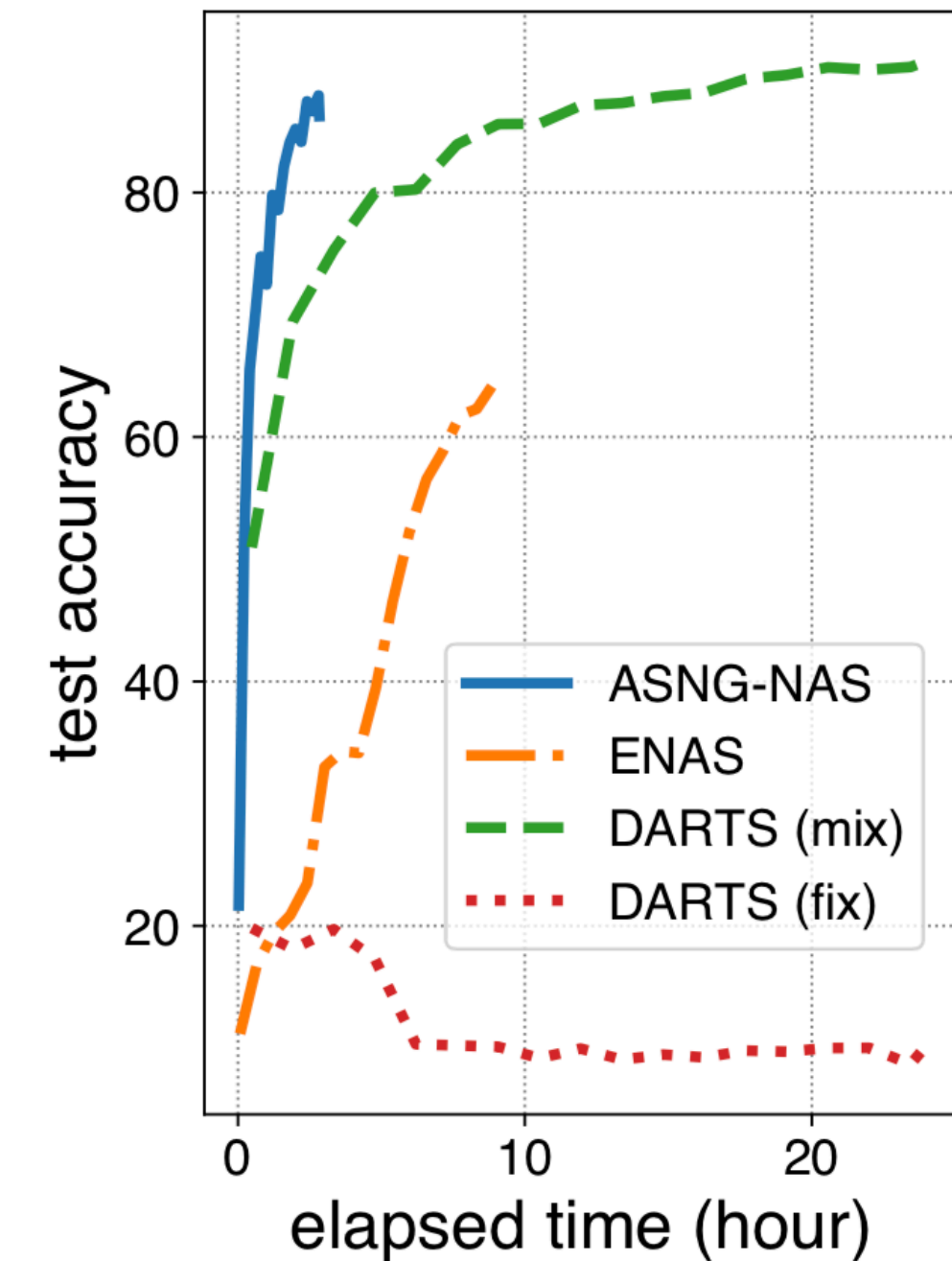


Figure 2: Transitions of test error against elapsed time in the architecture search phase.

The detail will be explained at Poster #53