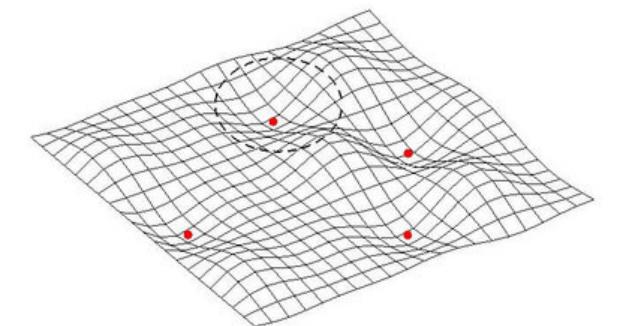
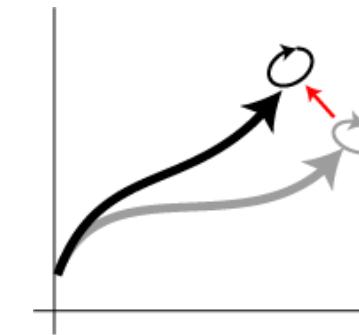
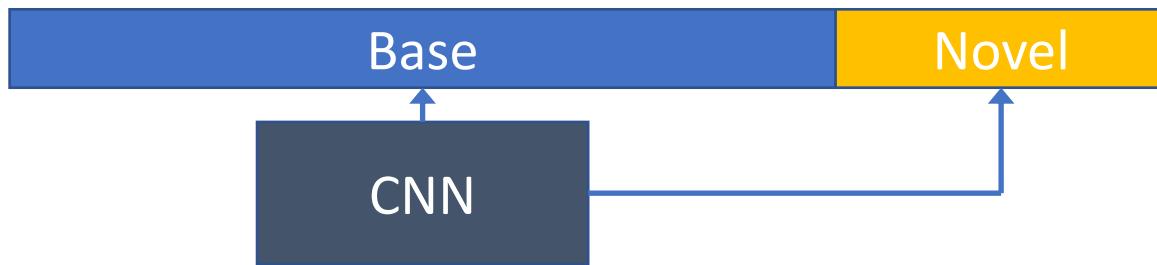


Incremental Few-Shot Learning with Attention Attractor Networks

Mengye Ren, Renjie Liao, Ethan Fetaya, Richard S. Zemel

University of Toronto and Vector Institute

- Testing **on only new classes** in “few-shot” is not natural.
- **Incremental few-shot learning:** learn new classes on top of old classes. **No access to the old data.**
- At each test episode, learn a linear classifier **until convergence**.
- Attention over base classes to form **attractor regularizers**.
- At the end of the episode, test on a **query set of both base and novel**.
- Use **recurrent backprop (RBP)** instead of **truncated BPTT** for learning more **stable loss functions**.
- **Learned regularizers** significantly reduce class interference.



Auto-Meta: Automated Gradient Based Meta Learner Search



Jaehong Kim¹, Sangyeul Lee¹, Sungwan Kim¹, Moonsu Cha¹, Jung Kwon Lee¹,
Youngduck Choi^{1,2}, Yongseok Choi¹, Dong-Yeon Cho¹, and Jiwon Kim¹



¹ SK T-Brain ² Yale University

Automated architecture
search

Gradient-based
Meta learning

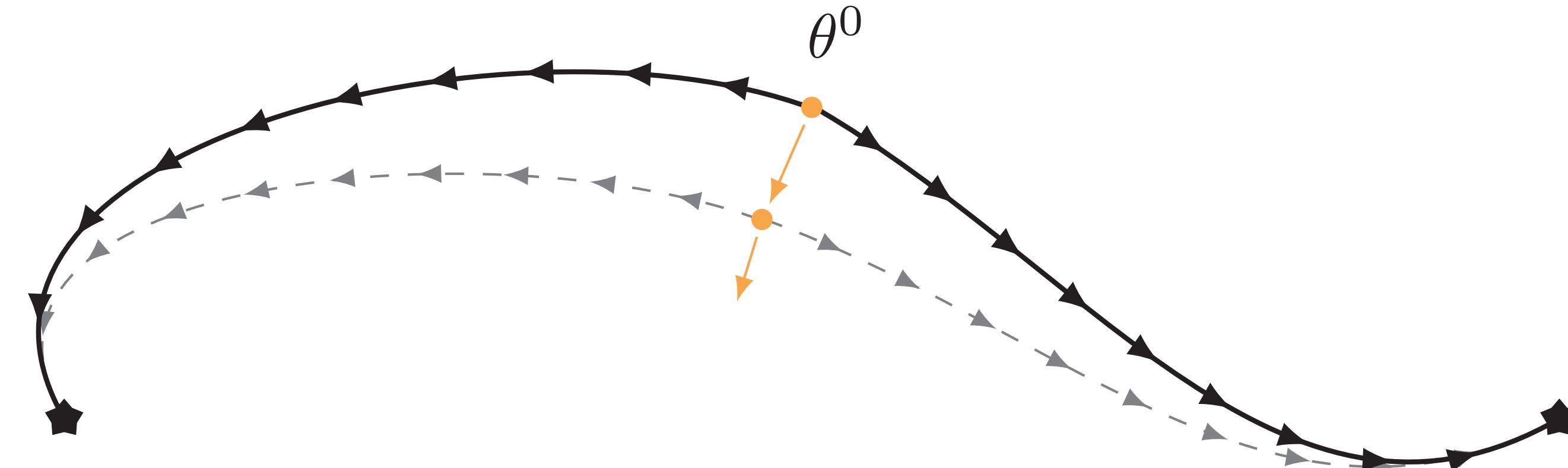
The diagram illustrates the methodology of Auto-Meta. Two arrows originate from the text boxes above and point downwards towards a central blue oval. The left arrow is green and points to the text 'Automated architecture search'. The right arrow is orange and points to the text 'Gradient-based Meta learning'. Both arrows converge on the central blue oval, which contains the text 'Performance improvement' and 'Few-shot image classification (Omniglot, Mini-ImageNet)'.

Performance improvement
Few-shot image classification
(Omniglot, Mini-ImageNet)

Transferring Knowledge across Learning Processes

Sebastian Flennerhag, Pablo G. Moreno, Neil D. Lawrence, Andreas Damianou

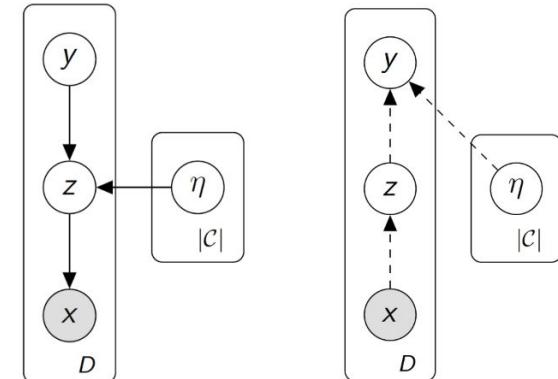
- We propose a framework for meta-learning across task geometries by learning from gradient trajectories
- We present *Leap*, a light-weight meta-learner that scales beyond few-shot learning to tasks requiring millions of gradient steps



Few-shot Learning For Free by Modelling Global Class Structure

Xuechen Li*, Will Grathwohl*, Eleni Triantafillou*, David Duvenaud, Richard Zemel

- Most approaches to few-shot classification use **episodic training**.
- We advocate for a simpler approach: a generative model over **all classes**: a VAE with a **mixture of Gaussians prior**.
- Few-shot learning is done by **variational inference**.
- Our model solves 3 tasks:
 - Few-shot classification
 - Few-shot generation
 - More realistic: **Few-shot integration**.
- Omniglot experiments:
 - On par with state-of-the-art on few-shot classification.
 - Largely outperform our baseline on few-shot integration.



TAEML: Task-Adaptive Ensemble of Meta-Learners

A I T R I C S

Workshop on Meta-Learning (MetaLearn2018)

Minseop Park / mike_seop@aitrics.com

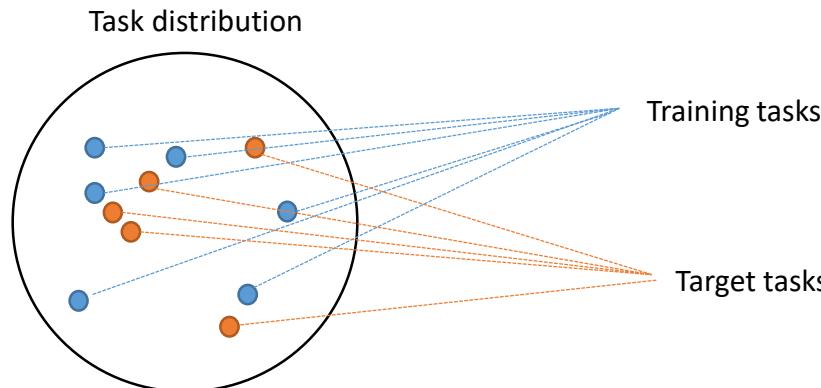


Fig1. Current meta-learning for few-shot classification

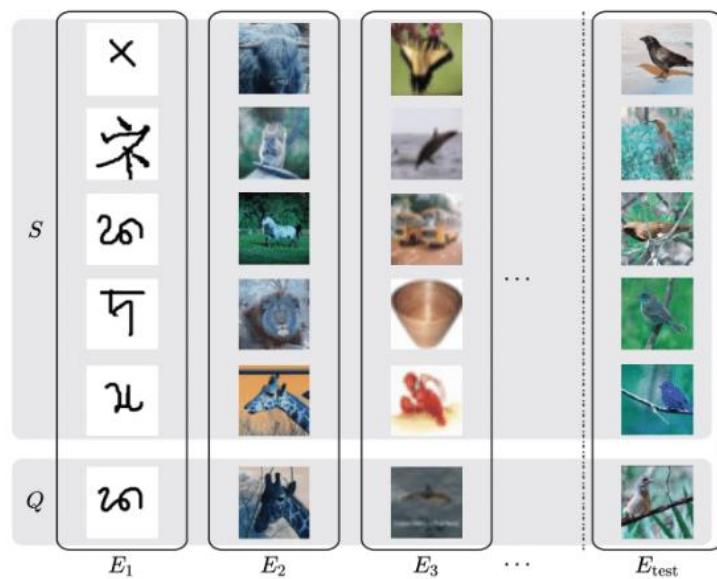


Fig2. Solving to few-shot classify the birds: Training all of the tasks won't be efficient

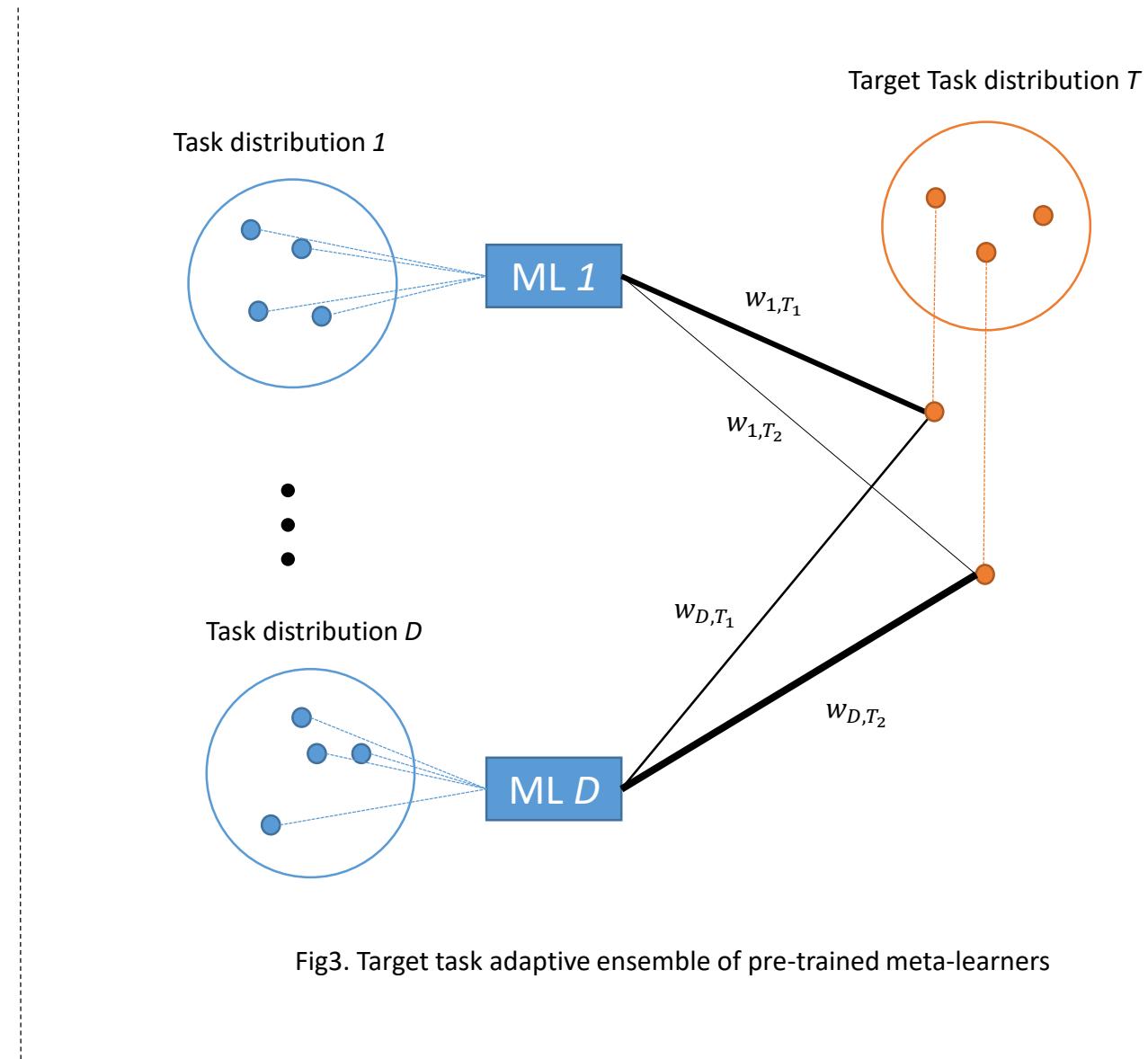


Fig3. Target task adaptive ensemble of pre-trained meta-learners

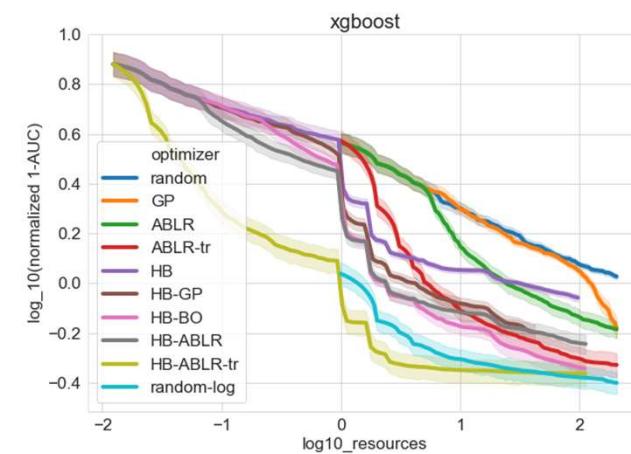
A Simple Transfer-Learning Extension of Hyperband

Lazar Valkov, Rodolphe Jenatton, Fela Winkelmolen, Cédric Archambeau

- Setting: Hyperparameter Optimisation
- Hyperband (HB):
 - Incrementally allocates more resources to the best-performing candidates initially taken from a pool of randomly sampled candidates.
 - Evaluates different number of initial candidates n_i for r_i
- We enhance HB with model-based sampling, using ABLR (Peronne *et al.*)

$$P(\mathbf{y}_t | \mathbf{w}_t, z, \beta_t) = \mathcal{N}(\Phi_{\mathbf{z}}(\mathbf{X}_t, \mathbf{r}_i) \mathbf{w}_t, \beta_t^{-1} I_{N_t}) P(\mathbf{w}_t | \alpha_t) = \mathcal{N}(\mathbf{0}, \alpha_t^{-1} I_D)$$

- Benefits:
 - Makes use of all data produced by a HB run
 - Can use data from past HB runs to learn better basis function
 - We don't use heuristics for low number of data points, nor to encourage exploration



Learned optimizers that outperform SGD on wallclock and test loss

Google AI



Luke Metz, Niru Maheswaranathan, Jeremy Nixon, C. Daniel Freeman, Jascha Sohl-Dickstein

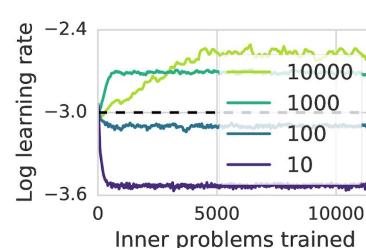
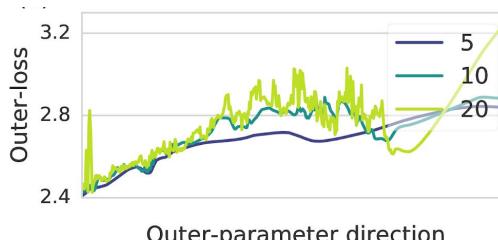
Existing optimizers are **hand designed**. Can we do better with **learning**?

One popular strategy for training such optimizers is to leverage gradients and **truncated backpropagation through time**.

These methods, however, are notoriously **unstable**!

Careful choice of step length is required:

- Long truncations: **exploding gradients**
- Short truncations: **biased gradients**

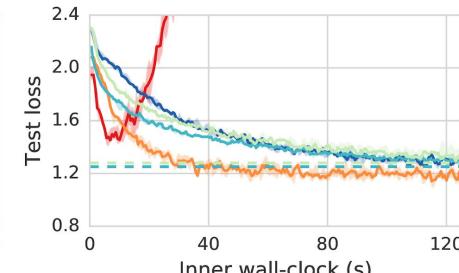
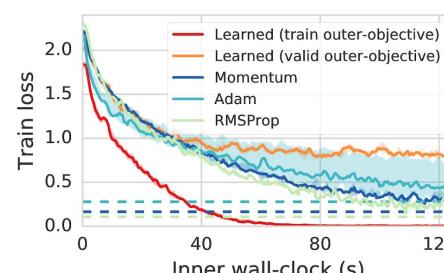


We use **variational optimization** to "smooth" the loss surface by convolving it with a Gaussian.

$$\mathcal{L}(\theta) = \mathbb{E}_{\tilde{\theta} \sim \mathcal{N}(\theta, \sigma^2 I)} [L(\tilde{\theta})]$$

To optimize this objective, we combine **multiple gradient estimators** with difference variances.

We train **simple** MLP-based learned optimizers that are **faster in wallclock time** and **generalize better** than existing hand-designed methods.



Define two gradient estimators:

- **reparameterization trick**
- **evolutionary strategies**

Combine them!

$$g_{\text{rp}} = \frac{1}{S} \sum_s \nabla_{\theta} L(\theta + \sigma n_s), \quad n_s \sim N(0, I)$$
$$g_{\text{es}} = \frac{1}{S} \sum_s L(\tilde{\theta}_s) \nabla_{\theta} \left[\log \left(N(\tilde{\theta}_s; \theta, \sigma^2 I) \right) \right], \quad \tilde{\theta}_s \sim N(\theta, \sigma^2 I)$$



Learning to Learn with Conditional Class Dependencies

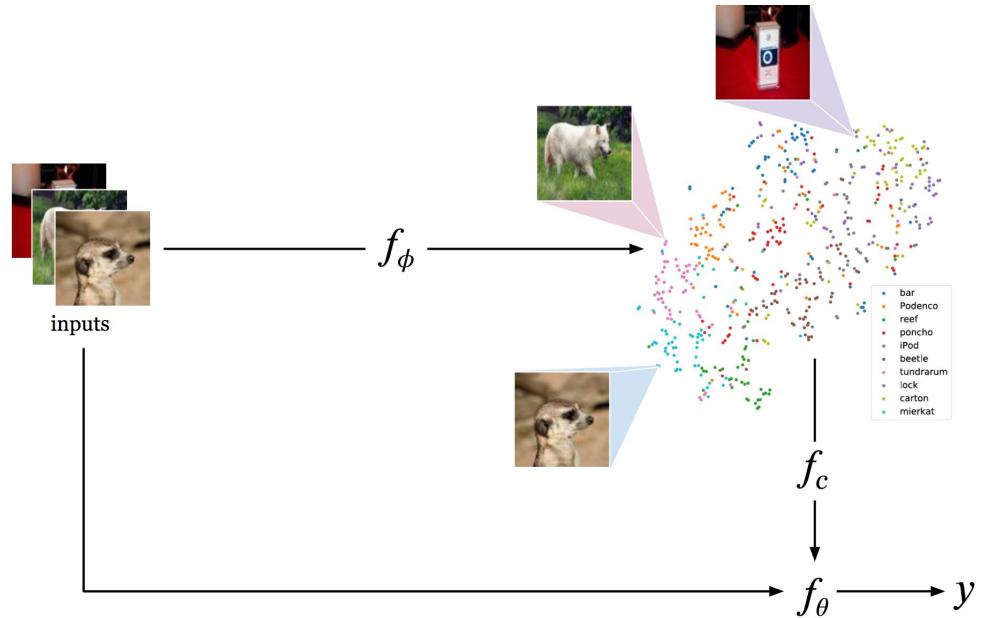
Xiang Jiang^{1,2}, Mohammad Havaei¹, Farshid Varno^{1,2}, Gabriel Chartrand¹, Nicolas Chapados¹, Stan Matwin²

¹Imagia Inc. ²Dalhousie University

Integrates **two views** of the data

The metric space captures **class dependencies**

Conditional batchnorm helps **class separation**





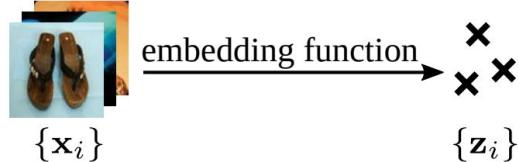
Unsupervised Learning via Meta-Learning

Kyle Hsu¹, Sergey Levine², Chelsea Finn²

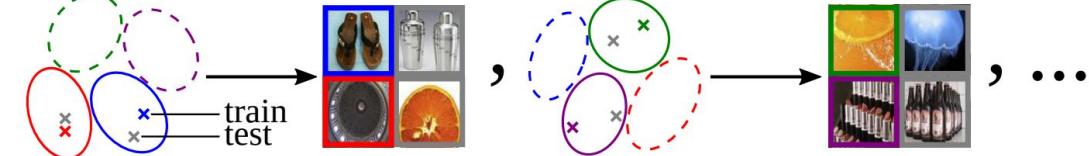
¹University of Toronto ²UC Berkeley

- Unsupervised learning is commonly used as pre-training for downstream learning.
 - We improve upon this by incorporating knowledge about the downstream task type: image classification.
- **Unsupervised meta-learning** via CACTUs: meta-learning over tasks constructed from unlabeled data.

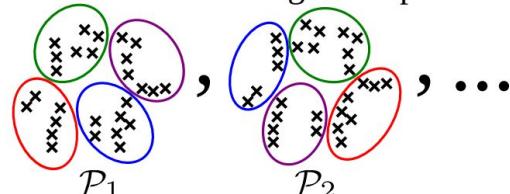
1. run embedding learning



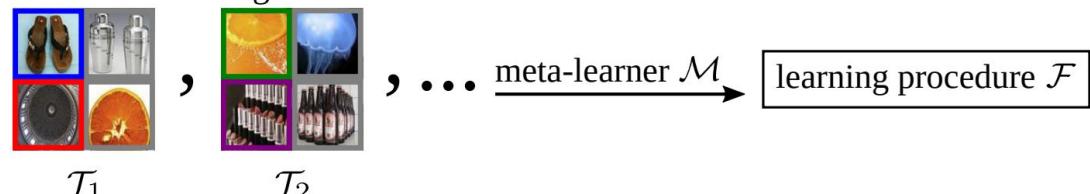
2b. automatically construct tasks without supervision



2a. cluster embeddings multiple times



3. run meta-learning on tasks



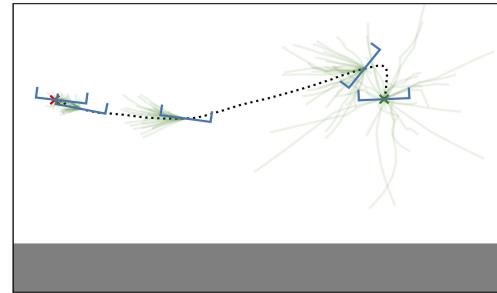
- Results: better than unsupervised learning, worse than supervised meta-learning

CAMeLiD: Control Adaptation via Meta-Learning Dynamics

James Harrison^{*,1}, Apoorva Sharma^{*,1},
Roberto Calandra², Marco Pavone¹



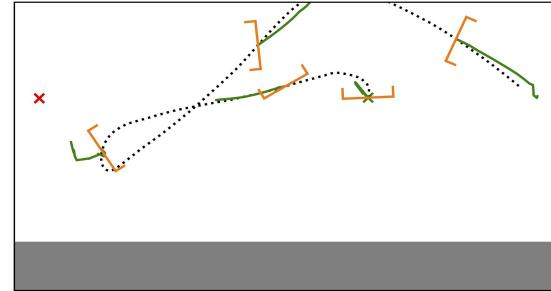
We develop a Bayesian meta-learning model that is capable of **fast, efficient online updates** and is trained for multi-step probabilistic predictions.



CAMeLiD controlling a quadrotor with a random attached mass. By incorporating model uncertainty into control, we successfully stabilize.

Using this model, we build a control algorithm that captures online model uncertainty and **automatically trades off safety and performance**.

Point estimate meta-learning-based control algorithm results in the quadrotor crashing.



Learning to Adapt in Dynamic, Real-World Environments Through Meta-Reinforcement Learning



Anusha Nagabandi*, Ignasi Clavera*, Simin Liu,
Ron S. Fearing, Pieter Abbeel, Sergey Levine, Chelsea Finn

Goal

Use **recent experiences** to quickly **adapt** to the current situation.

Train time: Learning to Adapt

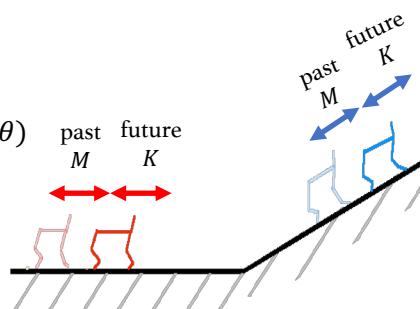
Meta-learn a dynamics model

Tasks: temporal windows

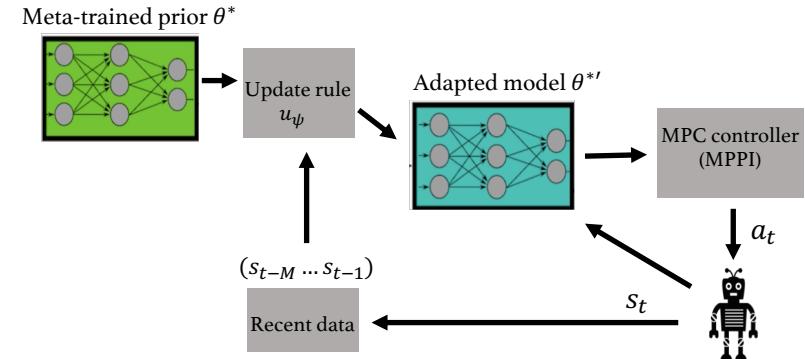
Objective:

$$\min_{\theta, \psi} E[L(D_T^{test}, \theta')] \text{ s.t. } \theta' = u_\psi(D_T^{tr}, \theta)$$

D_T^{test} → Future data
 D_T^{tr} → Past data



Test time: Meta-Model-Based RL



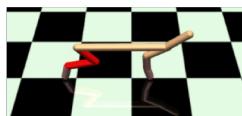
Experiments



Pier



Terrain slopes



Disabled



Crippled



Slope



Pose error



Payload



Missing leg

Learning to Design RNA

Frederic Runge*

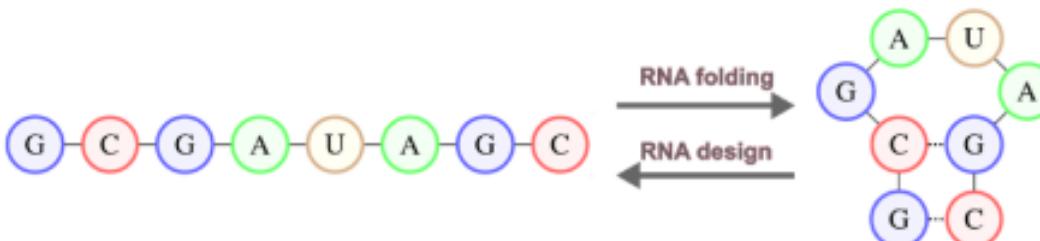
Danny Stoll*

Stefan Falkner

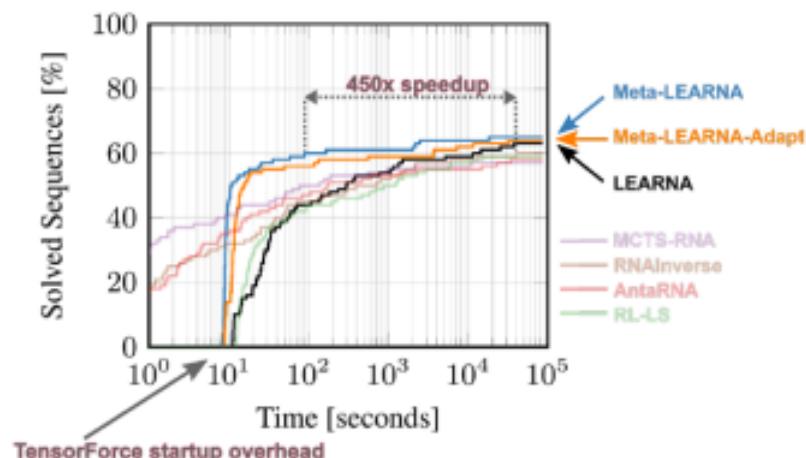
Frank Hutter



UNI
FREIBURG



- **Meta-learn** a policy across RNA Design tasks
- **AutoML** for joint optimization of:
 - Policy network architecture
 - RL formulation
 - Training Hyperparameters
- **New state-of-the-art** on three benchmarks

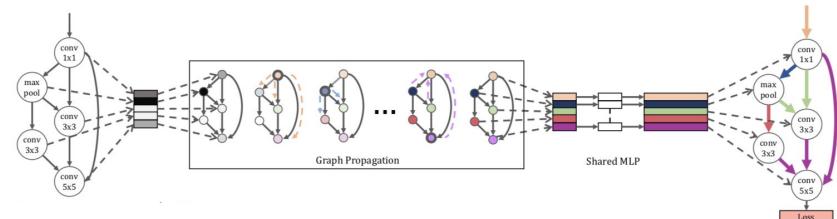


Graph HyperNetworks for Neural Architecture Search

Chris J. Zhang^{1,2}, Mengye Ren^{1,3}, Raquel Urtasun^{1,3}

¹ Uber Advanced Technologies Group ² University of Waterloo, ³ University of Toronto

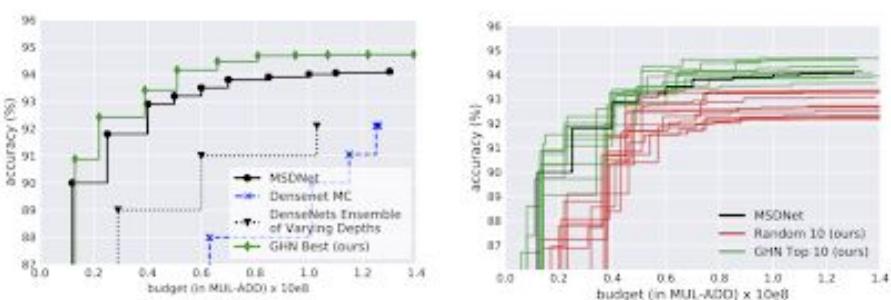
Graph HyperNetworks



Motivation:

- Neural architecture search is an expensive nested optimization
- Instead of using SGD to learn weights, use trained hypernetwork to generate weights
- Graph HyperNetworks (GHN) explicitly model the topology of architectures by learning on a computation graph representation

Anytime Prediction



NAS Benchmarks

CIFAR-10: Comparison with NAS methods which employ random search (top half) and advanced search methods (e.g. RL) (bottom half)

Method	Search Cost (GPU days)	Param $\times 10^6$	Accuracy
SMASHv1 (Brock et al., 2018)	?	4.6	94.5
SMASHv2 (Brock et al., 2018)	3	16.0	96.0
One-Shot Top (F=32) (Bender et al., 2018)	4	2.7 ± 0.3	95.5 ± 0.1
One-Shot Top (F=64) (Bender et al., 2018)	4	10.4 ± 1.0	95.9 ± 0.2
Random (F=32)	-	4.6 ± 0.6	94.6 ± 0.3
GHN Top (F=32)	0.42	5.1 ± 0.6	95.7 ± 0.1
NASNet-A (Zoph et al., 2018)	1800	3.3	97.35
ENAS Cell search (Pham et al., 2018)	0.45	4.6	97.11
DARTS (first order) (Liu et al., 2018b)	1.5	2.9	97.06
DARTS (second order) (Liu et al., 2018b)	4	3.4	97.17 ± 0.06
GHN Top-Best, 1K (F=32)	0.84	5.7	97.16 ± 0.07

ImageNet Mobile: Comparison with NAS methods which employ advanced search methods (e.g. RL)

Method	Search Cost (GPU days)	Param $\times 10^6$	FLOPs $\times 10^6$	Top 1	Top 5
NASNet-A (Zoph et al., 2018)	1800	5.3	564	74.0	91.6
NASNet-C (Zoph et al., 2018)	1800	4.9	558	72.5	91.0
AmoebaNet-A (Real et al., 2018)	3150	5.1	555	74.5	92.0
AmoebaNet-C (Real et al., 2018)	3150	6.4	570	75.7	92.4
PNAS (Liu et al., 2018a)	225	5.1	588	74.2	91.9
DARTS (second order) (Liu et al., 2018b)	4	4.9	595	73.1	91.0
GHN Top-Best, 1K	0.84	6.1	569	73.0	91.3

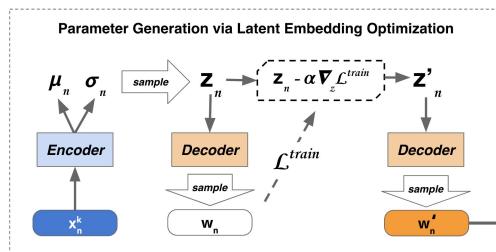
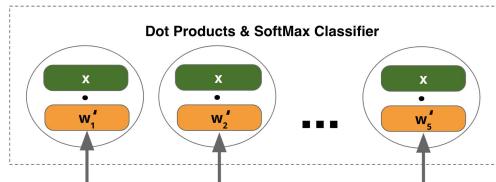
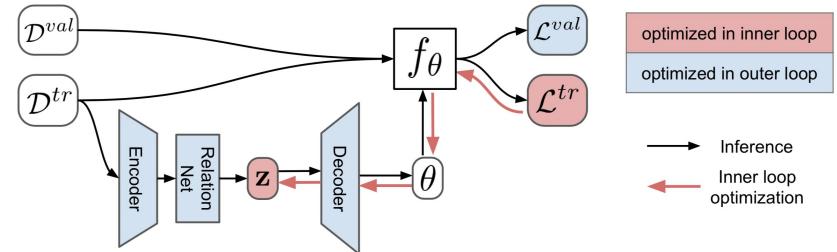
Meta-Learning with Latent Embedding Optimization (LEO)

Andrei A. Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, Raia Hadsell

We learn a data-dependent latent generative representation of model parameters, and perform gradient-based meta-learning in this low dimensional latent space.

The resulting approach, Latent Embedding Optimization (LEO), decouples the gradient-based adaptation procedure from the underlying high-dimensional space of model parameters.

LEO is *state-of-the-art* on both *minilmageNet* and *tieredImageNet* 5-way 1-shot and 5-shot classification tasks.



We are in the process of open-sourcing our embeddings and code!



DeepMind

Proximal Meta-Policy Optimization: ProMP

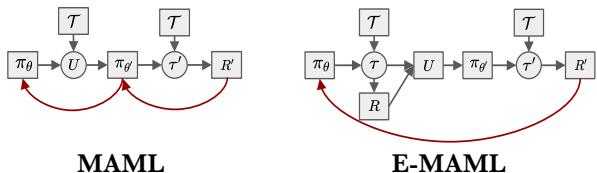
Jonas Rothfuss*, Dennis Lee*, Ignasi Clavera*,
Tamim Asfour, and Pieter Abbeel



Goal

- Analyze **credit assignment** in meta-reinforcement learning
- Develop a **new objective** that trains for the pre-update sampling distribution

Credit Assignent Sampling Distribution



Low Variance Curvature Estimator (LVC)

$$J_{\mathcal{T}}^{\text{LVC}}(\tau) = \sum_{t=0}^{H-1} \frac{\pi_\theta(a_t|s_t)}{\perp(\pi_\theta(a_t|s_t))} \left(\sum_{t'=t}^{H-1} r(s_{t'}, a_{t'}) \right) \quad \tau \sim P_{\mathcal{T}}(\tau)$$

- Meta-gradient with **low variance**
- Unbiased** closed to local optima

Proximal Meta-Policy Optimization: ProMP

ProMP Objective:

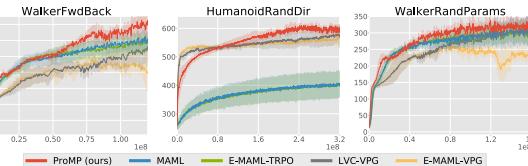
$$J_{\mathcal{T}}^{\text{ProMP}}(\theta) = J_{\mathcal{T}}^{\text{CLIP}}(\theta') - \eta \bar{D}_{KL}(\pi_{\theta_o}, \pi_\theta) \quad \text{s.t.} \quad \theta' = \theta + \alpha \nabla_\theta J_{\mathcal{T}}^{\text{LR}}(\theta)$$

Incorporates the benefits of:

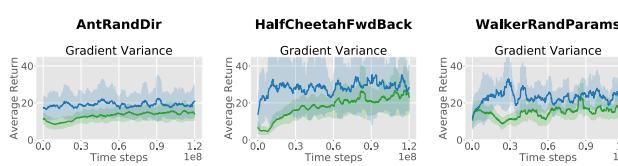
- Proximal Policy Optimization
- LVC Estimator

Experiments

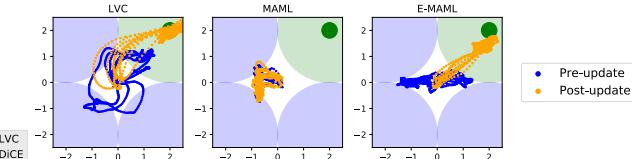
Performance Comparison



Variance Comparison



Exploration – Exploitation



Attentive Task-Agnostic Meta-Learning for Few-Shot Text Classification

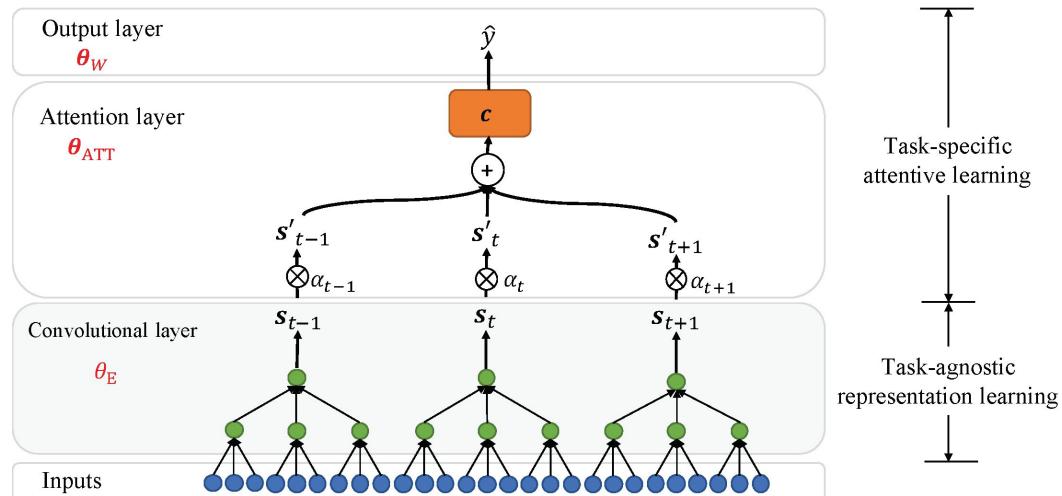


Xiang Jiang^{1,2}, Mohammad Havaei¹, Gabriel Chartrand¹, Hassan Chouaib¹, Thomas Vincent¹, Andrew Jesson,¹ Nicolas Chapados¹, Stan Matwin²
¹Imagia Inc. ²Dalhousie University

Task-agnostic representation learning

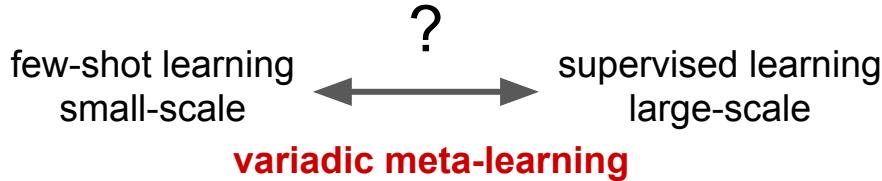
Task-specific attentive adaptation

Attention decouples the representation learning



Variadic Meta-Learning by Bayesian Nonparametric Deep Embedding

Kelsey Allen, Hanul Shin*, Evan Shelhamer*, Josh Tenenbaum

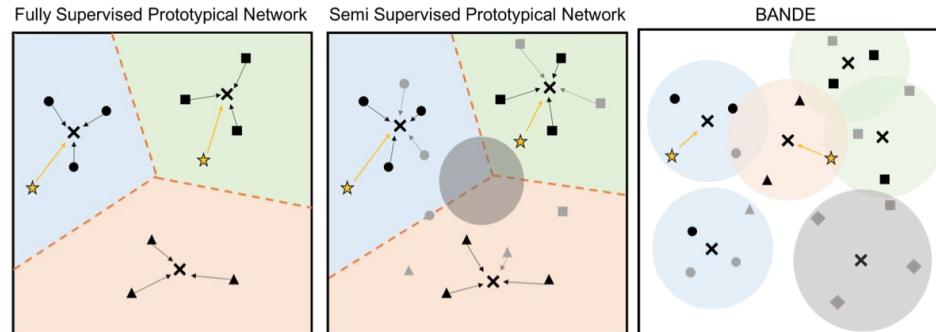


any-shot, any-way generalization
between meta-train and meta-test
with mixed supervision

experiments:

- from **5-way to 1692-way** and
from **1-shot to unsupervised** on Omniglot
- from **1-shot to 50-shot** on mini-ImageNet
- from **2-shot to 5000-shot** on CIFAR-10

with comparison of prototypes, MAML, graph nets,
and good old supervised learning



BANDE clusters labeled and unlabeled data into *multi-modal prototypes* that represent each class by a set of clusters instead of only one

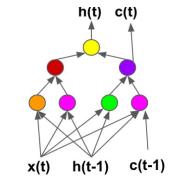
multi-modal prototypes
for alphabet and character recognition

Training	Testing	Proto. Nets	BANDE
Alphabet	Alphabet	64.9 ± 0.2	91.2 ± 0.1
Alphabet	Chars (20-way)	85.7 ± 0.2	95.3 ± 0.2
Chars	Chars (20-way)	94.9 ± 0.2	95.1 ± 0.1

From Nodes to Networks: Evolving Recurrent Neural Networks

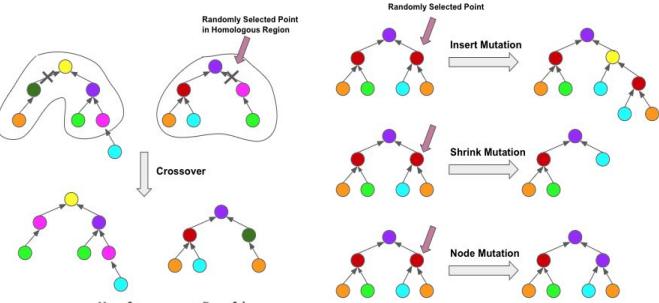
Aditya Rawal*, Risto Miikkulainen*
 aditya.rawal@uber.com, risto@cs.utexas.edu

* Work done at Sentient Technologies



Recurrent Cell
as Tree

Evolve

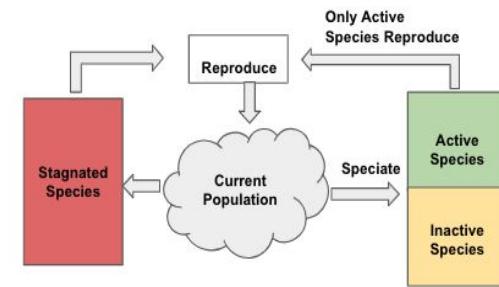


Crossover

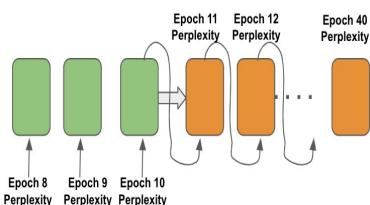
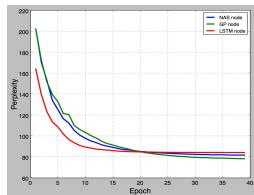
Mutation

$$\delta_T(T_i, T_j) = \beta \frac{N_{i,j} - 2n_{S_{i,j}}}{N_{i,j} - 2} + (1 - \beta) \frac{D_{i,j} - 2d_{S_{i,j}}}{D_{i,j} - 2}$$

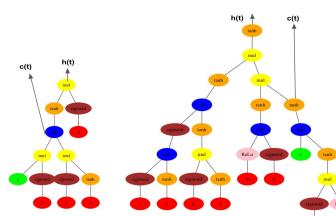
(Tajuddin et al., 2015)



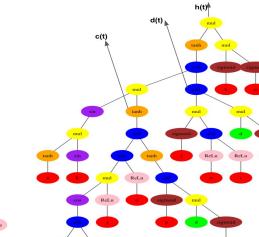
Encourage Search for Novel Cells



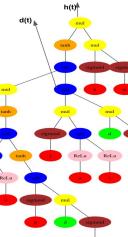
Meta-LSTM: Seq2Seq model to predict learning curve.
Speeds-up search by **4X**.



LSTM



NAS Cell



Evolved
Cell

Language Modeling

Transfer to Music



Music