AutoRL Tutorial @ AutoML

Aleksandra Faust Senior Staff Research Scientist <u>faust@google.com</u> <u>https://www.afaust.info/</u>

Google Brain, July 2022

Google Research

Yingjie Miao Senior Software Engineer <u>vingjiemiao@google.com</u> Richard Song Research Scientist xingyousong@google.com https://xingyousong.github.io/

Content

Why AutoRL?

- Why does RL need automation?
- AutoRL problem formulation

AutoRL training methods and tools

What parts of RL we can automate?

What are open problems?

TI;Dr;



Why AutoRL?





Learn through trial and error by interacting with the environment.

RL Primer - What is RL?



Image credit: Waymo



Image credit: UVM



fetch



Image credit: Loon

What is Reinforcement Learning?



Learns a policy that maximizes expected cumulative return.

Components of RL

Value functions / estimators:

- State value, V: the potential payoff for the state
- State action value, Q: potential payoff for the action

Exploration vs. Exploitation Trade off

Choice of what to learn

- Model model-based RL
- Value functions value iteration RL
- Policy policy search, or policy iteration
- Value and policy actor critic methods.





Image credit: Open Al

How is RL different from SL? And why is it difficult to train?



Closed loop training and distributional shift in data



Every training iteration gathers the data from changed environment.

Closed loop training and distributional shift in data



More complicated RL algorithm.

Watch the video at for an example of distributional shift at the evaluation time:

https://www.youtube.com/watch?v=xN-OWX5gKvQ&t=1s



Closed loop evaluation and precision / latency trade off.



Environment distribution shift grows with the policy latency.

Bigger not always better. Find the sweet spot.

Function approximations and PoMDP $(S, A, P, R, O, \Omega, \rho_o, T, \gamma)$



PoMDP doesn't say anything about function approximations.

Function approximations and PoMDP $(S, A, P, R, O, \Omega, \rho_o, T, \gamma)$



PoMDP doesn't say anything about function approximations.

Experimentation with NN architectures.

Function approximations and PoMDP $(S, A, P, R, O, \Omega, \rho_o, T, \gamma)$



Simulator design and Markovian property violation.

More complicated stack.

Application Implications: Safety and Generalization



Often can cause damage, and needs to work in unstructured environments.

Training environment design and generalization methods.

How is RL different from SL? And why is it difficult to train?



Reality of Training and Evaluating RL Agents



- RL-focused components (Algorithm, NAS)
- Task definition (reward)
- Agent (observations, simulator, actions)
- Environment (curriculum)

Reality of Training and Evaluating RL Agents



- RL-focused components (Algorithm, NAS)
- Task definition (reward)
- Agent (observations, simulator, actions)
- Environment (curriculum)

Reality of Training and Evaluating RL Agents



- RL-focused components (Algorithm, NAS)
- Task definition (reward)
- Agent (observations, simulator, actions)
- Environment (curriculum)

$$\max_{\zeta} f(\zeta, \theta^*) \quad \text{s.t.} \quad \theta^* \in \argmax_{\theta} J(\theta; \zeta)$$

AutoRL Problem Formulation



Hyperparameter Taxonomy by Type

Learning Hyperparameters Low-parameter regime

High-parameter regime



Image credit: Colton Bishop

Image credit: [Co-Reyes et al., 2021]

Hyperparameter Taxonomy by Training-step Dependence

Is there an optimal hyperparameter for the entire inner training loop?



AutoRL Cost

High-parameter regime

 $\theta: S \rightarrow List[R]$

SelectList

at

NN: $S \rightarrow List[R]$

St

class **stable_baselines.sac.SAC**(policy, env, gamma=0.99, learning_rate=0.0003, buffer_size=50000, learning_starts=100, train_freq=1, batch_size=64, tau=0.005, ent_coef='auto', target_update_interval=1, gradient_steps=1, target_entropy='auto', action_noise=None, random_exploration=0.0, verbose=0, tensorboard_log=None, _init_setup_model=True, policy_kwargs=None, full_tensorboard_log=False, seed=None, n_cpu_tf_sess=None) [source]

class stable_baselines.common.policies.BasePolicy(sess, ob_space, ac_space, n_env, n_steps, n_batch, reuse=False, scale=False, obs_phs=None, add_action_ph=False) [source]

28 hyperparameters w/o NAS, algorithm, and env design



AutoRL Design Considerations



Hyperparameters -- everything not in PoMDP.

- Categorical, symbolic, continuous, differentiable
- Stationary or nonstationary
- Co-dependencies.

Expensive evaluations.

- Search in large spaces
- Long evaluations (hours or days)

AutoRL Training Methods



Training Methods

- Search
- Gradient-free methods
 - Evolutionary methods
- Optimization methods
 - Bayesian Optimization
 - Gradient-based (MAML)
- Multi Arm bandits
- Population-based Training



Training Methods - Search



In multidimensional search spaces, random search provides more information about the objective function.

Image credit: [Parker-Holder et al., JAIR 2022]



Training Methods - Evolutionary Methods



Training Methods - Multi-Arm Bandits



Stateless action selection

Google Research

Zhou, 2015

Training Methods - Bayesian Optimization

Regressor (e.g. Gaussian Process): Uncertainty estimates of the objective function.

Acquisition function: Explore / exploit trade-off on f.

Next sample recommendation: Argmax of the acquisition function.



Training Methods - Optimization - MAML

Algorithm 1 Model-Agnostic Meta-Learning

Require: $p(\mathcal{T})$: distribution over tasks

Require: α , β : step size hyperparameters

- 1: randomly initialize θ
- 2: while not done do
- 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
- 4: for all \mathcal{T}_i do
- 5: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ with respect to K examples
- 6: Compute adapted parameters with gradient descent: $\theta'_i = \theta \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- 7: end for
- 8: Update $\theta \leftarrow \theta \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
- 9: end while



Figure 1. Diagram of our model-agnostic meta-learning algorithm (MAML), which optimizes for a representation θ that can quickly adapt to new tasks.

Credit: [Finn et al., ICML 2017]

Training Methods - Population-based Methods



PBT jointly optimizes a population of models and their Hyperparameters. Nonstationary parameters. [Jaderberg et al, 2017]

Training Methods - Trade-offs

Class	Algorithm properties				ies	What is automated?
Random/Grid Search (4.1)	የ የተ		\Rightarrow	1	÷	hyperparameters, architecture, algorithm
Bayesian Optimization (4.2)	የ የተ		\Rightarrow	1	÷	hyperparameters, architecture, algorithm
Evolutionary Approaches (4.3)	የ የ የ		\Rightarrow	1	\approx	hyperparameters, architecture, algorithm
Meta-Gradients (4.4)	f	∇	\rightarrow	•	\approx	hyperparameters
Blackbox Online Tuning (4.5)	ť		\rightarrow	1	\approx	hyperparameters
Learning Algorithms (4.6)	የ የየ		\Rightarrow	•	÷	algorithm
Environment Design (4.7)	የ የየ		\Rightarrow	•	\approx	environment

† only uses a single trial, *†††* requires multiple trials

- ∇ requires differentiable variables, \blacksquare works with non-differentiable hyperparameters
- \rightrightarrows parallelizable \rightarrow not parallelizable
- \checkmark works for any RL algorithm, \bigcirc works for only some classes of RL algorithms
- \doteq static optimization, \approx dynamic optimization

What to learn?


What to learn?

RL-focused components

- Tools for Hyperparameters
- Neural Architecture Search (NAS)
- RL Algorithm Learning

Task components: Reward

Agent-based environment learning

Environment-based learning

Performance

Auto RL



Environment

Observations, Rewards

Hyperparameter Tuning Packages

- Services: Host algorithms on a server.
 - More flexible + scalable
 - Additional engineering complexity
- **Frameworks:** Execute entire optimization (both algorithm + user evaluation)
 - Convenient, full automation
 - Requires knowledge of entire evaluation pipeline
- <u>Libraries:</u> Implement blackbox optimization algorithms
 - Offer most freedom
 - Lack scalability features / limited to single machine











Thou shall learn hyperparameters.

Open Source Vizier Service

OSS Vizier Service: Allows distributed multi-client + reliable tuning Supports most types of blackbox optimization





Thou shall learn hyperparameters.



Demo w/ OSS Vizier + Dopamine RL

(More in github.com/google/vizier/tree/main/demos)





Atari100K

Client + IDs

Evaluator +

Search Space +

Algorithm

Evaluation Loop

Thou shall learn hyperparameters.

Image Credit: [Machado et al, JAIR 2018]

What to Learn? Neural Network Architectures



What to learn? Neural Network Architectures

Architectures matter for:

- Inference Speed
- Generalization
- Sample Complexity
- Adaptation



Credit: [Yu et al. CoRL 2019]



Credit: [Ghost Robotics]



Credit: [Cobbe et al. ICML 2020]

Efficient Methods in NAS for RL

Original trial-by-trial formulation: Way too slow + expensive.

- CIFAR-10 already takes 2000+ GPU days to do NAS! [Zoph et al, CVPR 2018]
- In RL, noisy evaluation only makes it worse!

What are some efficient methods?



Credit: Depositphotos (Stock)



Low Parameter Regime

Low Parameter Regime (~10K weights): Optimize architecture + weights simultaneously

- Evolutionary Methods
 - NEAT [Stanley + Miikkulainen, EC 2002]
 - Regularized Evolution [Real et al., AAAI 2019]
 - ES-ENAS [Song et al, 2022]

Environments



5 2 2 Cart Position Long Pole Angle Short Pole Angle Bias

Stanley + Miikkulainen, Evolutionary Computation 2002



Architectures

Gaier + Ha, NeurIPS 2019

Caveat: Zero-order methods suffer under too many weight parameters.

High Parameter Regime: DARTS

High Parameter Regime (1M+ weights): Gradients for weight-training

 DARTS [Liu et al, ICLR 2019] used for RL in [Miao et al, AutoML-Conf 2022]: Minimally invasive!





Credit: [Cobbe et al. ICML 2020]



Credit: [Miao et al. AutoML-Conf 2022]

#feature_network = base_policies.make_impala_cnn_network(depths=[16, 32, 32])
feature_network = darts_policies.DartsIMPALACNN(depths=[16, 32, 32])

github.com/google/brain_autorl/tree/main/rl_darts

High Parameter Regime (Open Questions)

Possible Methods for NAS in RL? Little work so far.

- Blackbox Optimization? How to deal w/ strong noise?
 - Regularized Evolution [Real et al, AAAI 2019]?
 - BayesOpt?
- Other weight sharing methods? How to make simple?
 - ENAS [Pham et al, ICML 2018]?
- Others? How to scale to modern DL?
 - HyperNEAT [Stanley et al, Artificial Life 2009]?





Credit: CMU ML Blog, 2020



Credit: Roos, Blogpost on HyperNEAT, 2020

What to Learn? RL Loss Functions

Performance

Main idea: RL loss function, $\mathcal{L}(\theta; \zeta)$ is a parameterized object.



Search space millions of parameters and ~10¹⁰⁰ programs [Garau-Luis et al., 2022]

Careful design + pruning + optimization.



Loss function adaptive to environment and agent history.





MetaGenRL: population of agents train a single meta-objective [Kirsch et al., 2020]

Require Jacobian vector product computation. Readily available.

Some cross domain generalization.



Learned Policy Gradient: what and how to predict [Oh et al., 2020]



DSL for curiosity module trained w/ regressor. Search space combines NN, buffers, custom loss,

L2 Distance

etc. [Alet et al. 2020]

Interpretable loss functions.

target

Action Prediction Loss

Add To Loss

First time large databases of loss functions. github.com/jcoreves/evolvingrl

Loss function as a DAG. Lots of tricks due to size of

• Mutated

Hurdle Env

WWW

Lunit have

Spawn new training agent

Population of RL

Algorithms {L



Main idea: View RL algorithms (e.g. DQN) as points in a hyperparameter space. To discover algorithms is to do HP tuning in that space.



Main idea: View RL algorithms (e.g. DQN) as points in a hyperparameter space. To discover algorithms is to do HP tuning in that space.

Result: Co-Reyes et al. (2021) discovered new value-based RL algorithms that outperform baselines. These algorithms are interpretable, transferable, and implementable.



Main idea: View RL algorithms (e.g. DQN) as points in a hyperparameter space. To discover algorithms is to do HP tuning in that space.

Result: Co-Reyes et al. (2021) discovered new value-based RL algorithms that outperform baselines. These algorithms are interpretable, transferable, and implementable.

In this Tutorial:

- Algorithm representation and search space design
- Optimization method and tricks
- Code and colab example



Algorithm representation and the search space

Algorithm (loss function) as a DAG



$$L_{DQN} = (Q_{\theta}(s_t, a_t) - (r_t + \gamma * \max_a Q_{\theta'}(s_{t+1}, a)))^2$$

Algorithm representation and the search space

Algorithm (loss function) as a DAG



How to turn this into a HP and vice versa?

Each DAG can be viewed as a sequence of categorical HPs, one for each intermediate node (and the output node).



 $L_{DQN} = (Q_{\theta}(s_t, a_t) - (r_t + \gamma * \max_a Q_{\theta'}(s_{t+1}, a)))^2$

Algorithm representation and the search space

Algorithm (loss function) as a DAG



$$L_{DQN} = (Q_{\theta}(s_t, a_t) - (r_t + \gamma * \max_{a} Q_{\theta'}(s_{t+1}, a)))^2$$

How to turn this into a HP and vice versa?

Each DAG can be viewed as a sequence of categorical HPs, one for each intermediate node (and the output node).



For each intermediate node: which op to use (+, -, ...)? Which *previous* nodes as inputs?

For the output node: which intermediate node to output?

Encoding-decoding logic is necessary for mapping between HPs and DAGs.

Search space detail

Operation	Input Types	Output Type	
Add	Ж, Ж	X	
Subtract	Ж, Ж	X	
Max	Ж, Ж	X	
Min	Ж, Ж	X	
DotProduct	Ж, Ж	\mathbb{R}	
Div	Ж, Ж	X	
L2Distance	Ж, Ж	\mathbb{R}	
MaxList	$List[\mathbb{R}]$	\mathbb{R}	
MinList	$List[\mathbb{R}]$	\mathbb{R}	
ArgMaxList	$List[\mathbb{R}]$	Z	
SelectList	$List[X], \mathbb{Z}$	X	
MeanList	List[X]	X	
VarianceList	List[X]	X	
Log	X	X	
Exp	X	X	
Abs	X	X	
(C)NN: $\mathbb{S} \to List[\mathbb{R}]$	S	$List[\mathbb{R}]$	
$(\mathrm{C})\mathrm{NN}{:}\mathbb{S}\to\mathbb{R}$	S	\mathbb{R}	
$(\mathrm{C})\mathrm{NN}:\mathbb{S}\to\mathbb{V}$	V	\mathbb{V}	
Softmax	$List[\mathbb{R}]$	P	
KLDiv	\mathbb{P},\mathbb{P}	\mathbb{R}	
Entropy	\mathbb{P}	R	
Constant		1, 0.5, 0.2, 0.1, 0.01	
MultiplyTenth	X	X	
Normal(0, 1)		R	
Uniform(0, 1)		R	

Operations:

- Covers many known value-based RL algorithms.
- Typed inputs/outputs help to prune the search space.
- In practice, need to set a max size of the DAG (~20 worked).

This completes the specification of the HP search space. Next:

- Define an optimization objective.
- Plug in your favorite HP tuning algorithm!

Training overview



Optimization objective: sum of normalized returns on meta-training environments

Optimization method: Regularized Evolution

Distributed training is paramount (~300 CPUs for 3 days, ~20K candidates)

Speedup tricks:

- Hurdle environment (cartpole) fail fast!
- Functional equivalence hashing and caching dedup!

Discovered Algorithms

DQNReg: regular DQN + a regularization term

 $L_{\rm DQNReg} = 0.1 * Q(s_t, a_t) + \delta^2$

DQNReg generalizes to unseen classic Control, MiniGrids and Atari tasks.

Trained on non-image environments. Not tuned to Atari games.

Env	DQN	DDQN	PPO	DQNReg
Asteroid	1364.5	734.7	2097.5	2390.4
Bowling	50.4	68.1	40.1	80.5
Boxing	88.0	91.6	94.6	100.0
RoadRunner	39544.0	44127.0	35466.0	65516.0





DQNReg

Google Research

Baselines taken from reported numbers.



In code and colab





1. Search space definition



97 for i in range(search_program_length): 98 program_lst.append(99 pg.oneof(PyGlove allows arbitrarily complex HP 100 product_input_ops(inputs, existing_ops, i, operators, freeze_ops)))

A list of nodes, each node is a **tunable object**, as indicated by **pg.oneof**(**<list>**)

2. Search algorithm (Regularized Evolution)

282	<pre>return pg.evolution.Evolution(</pre>
283	reproduction=(
284	# Tournament selection and mutation.
285	<pre>pg.evolution.selectors.Random(tournament_size, seed=seed) >></pre>
286	<pre>pg.evolution.selectors.Top(1) >> graph_mutator),</pre>
287	<pre>population_init=(graph_generator, population_size),</pre>
288	population_update=(
289	# Pop out oldest individual and update functional equivalence cache.
290	<pre>pg.evolution.selectors.Last(population_size) >> update_cache))</pre>

In code and colab (continued)



End-to-end colab:

https://github.com/google/brain_autorl/blob/main/evolving_rl/EvolvingRL_Demo.ipynb



Careful design + pruning + optimization.

Generalize across domains in discrete action spaces. PG and Actor/Critic more challenging. [Garau-Luis et al., 2022]

What to Learn? Rewards

Performance



Main idea: Learn intrinsic rewards that maximize the excertic rewards.



Reward tuning is difficult.

What to learn? Task Components. Rewards.



[Zhen et al. 2018] adds intrinsic reward neural network. Inner and outer: GD.

S

[Chiang et al. 2019] Feature-based intrinsic reward w/ ES.

BiPaRS [<u>Hu et al. 2020</u>]

Feature-based reward shaping function w/ meta-learning.

7.5

Reward shaping helps for more complex tasks. Exact method is less important.

Learn the intrinsic reward that completes task.



Single trial: Equivalent of 12 days of training. 12 hours wall time.1000 trials: Equivalent of 32 years of training experience. 5 days wall time.

[Learning Navigation Behaviors End to End with AutoRL, Chiang*, Faust,* Fiser, Francis, RA-L/ICRA 2019] [Preprint, Video] Google Research



Zero-shot transfer





Generalizes to real world, unseen environment. Adapts to changes on the go.



Differential drive.

[Learning Navigation Behaviors End to End with AutoRL, Chiang*, Faust,* Fiser, Francis, RA-L/ICRA 2019 Preprint, Video]

Kinodynamic constraints.

Works on variety of robots.

How general is intrinsic reward learning?

SAC [Haarnoja et al, 2018] PPO [Schulman et al, 2017]

DDPG [Lillicrap et al, 2015]



Intrinsic reward helps in more complex tasks.

Learn the rewards instead of hyper-parameter tuning, on a fixed training budget.

Codependency between RL algorithm and intrinsic reward.

Simple objective performs almost as well -- without hand engineering.

[Evolving Rewards to Automate Reinforcement Learning, Faust, Francis, Mehta, 6th AutoML @ ICML 2019] [Preprint, Video] Google Research



What to Learn? Environment-based learning.



What to learn? Agent-based environment learning.



Dynamics

Abstract

We explore building generative neural networks models of popular neinforcement learning environments. Our world model can be trained quickly in an unsupervised madel can be trained of the environment. By using features extracted from the world model as inputs to an agent, extracted from the world model as inputs to an agent, extracted and the second second second second second second from the world model as the second second second from the world model and the second second second from the world model and transfer this second second from the world model and transfer this poly back into the actual environment.

An interactive version of this paper is available at

Figure 1. A World Model, from Scott McCloud's Understanding

Comics (McCloud, 1993; E. 2012)

World Models learn vision, memory, and environment w/ VAE. [<u>Ha and Schmidhuber, 2018</u>]

What to Learn? Curriculum.

Order of tasks and environments for the agent to train on.







Combinatorial Multi-Objective Evolutionary Algorithm (CMOEA) [<u>Huizinga and Clune, 2019</u>] [Kanitscheider et al, 2019] Learning progress based multi-task curriculum PAIRED constrained environment generation and RL agent. [Dennis et al, 2020]

Improves generalization. Sensitive to learning methods.
What to Learn? Combination of components.





LEARNA: reward + neural architecture + learning hyperparameters [Runge et al, 2019]

Auto-pilot: NAS + accelerator [Krishnan et al, 2022]

Summary of What to Learn?

What we know

Do learn learning HP.

Do learn intrinsic rewards for complex problems. Codependency with RL algorithm.

Do use curriculum for complex problems. Sensitive to task.

Loss function learning has promising generalization results in discrete action spaces. Requires careful design + pruning + optimization.

What we don't know

Very little work in NAS.

Very little work in agent learning.

What are important components to learn?

Which HP are nonstationary?

How to learn many components?

What are the best learning methods for each component?

Open Problems



Open Problems

- Theoretical understanding
 - HP co-dependencies.
 - Learning method appropriateness.
- Practical benchmarks
 - Non-trivial RL problem, and tractible outer loop.
 - Metrics.
- Nurture vs. nature trade off
- Integrative methods
 - Population-based training, multi-generational evolution + GD
 - Larger groups of HP, rewards, models, full algorithms.

Summary



TL;Dr;

Fundamental difference between RL and SL - more knobs and interdependencies. Stability of training and generalization due to non-stationarity and out of distribution.

Computational cost (AutoRL) = cost (AutoML) x cost (RL)

Hyperparameter taxonomy: numerical, symbolic, neural. Stationary or nonstationary.

Learn selected parts of the system: learning HP, NAS, loss, rewards, agent and environment components.

Future: Co-dependencies. Nonstationary hyperparameters. Scaling up.

Demos: github.com/google/vizier/tree/main/demos github.com/google/brain_autorl/blob/main/evolving_rl/EvolvingRL_Demo.ipynb github.com/jcoreyes/evolvingrl



Ten+ years of full AutoRL training. Joint curriculum and nonstationary hyperparameter training.

Acknowledgements



Yingjie Miao



Richard Song



John D. Co-Reyes



Jack Parker-Holder

Automated Reinforcement Learning (AutoRL): A Survey and Open Problems

Jack Parker-Holder University of Oxford

Raghu Rajan University of Freiburg

Xingyou Song Google Research, Brain Team

André Biedenkapp University of Freiburg

Yingjie Miao Google Research, Brain Team

Theresa Eimer Leibniz University Hannover

Baohe Zhang University of Freiburg

Vu Nguyen Amazon Australia

Roberto Calandra Meta AI

Aleksandra Faust Google Research, Brain Team

Frank Hutter University of Freiburg & Bosch Center for Artificial Intelligence

Marius Lindauer Leibniz University Hannover JACKPH@ROBOTS.OX.AC.UK

RAJANR@CS.UNI-FREIBURG.DE

XINGYOUSONG@GOOGLE.COM

BIEDENKA@CS.UNI-FREIBURG.DE

YINGJIEMIAO@GOOGLE.COM

EIMER@TNT.UNI-HANNOVER.DE

ZHANGB@CS.UNI-FREIBURG.DE

VUTNGN@AMAZON.COM

RCALANDRA@FB.COM

SANDRAFAUST@GOOGLE.COM

FH@CS.UNI-FREIBURG.DE

LINDAUER@TNT.UNI-HANNOVER.DE

Abstract

The combination of Reinforcement Learning (RL) with deep learning has led to a series of impressive feats, with many believing (deep) RL provides a path towards generally capable agents. However, the success of RL agents is often highly sensitive to design choices in the training process, which may require tedious and error-prone manual tuning. This makes it challenging to use RL for new problems and also limits its full potential. In many other areas of machine learning, AutoML has shown that it is possible to automate such design choices, and AutoML has also yielded promising initial results when applied to RL. However, Automated Reinforcement Learning (AutoRL) involves not only standard applications of AutoML but also includes additional challenges unique to RL, that naturally produce a different set of methods. As such, AutoRL has been emerging as an important area of different set of methods.

©2022 AI Access Foundation. All rights reserved.



Thank You

Fundamental difference between RL and SL - more knobs and interdependencies. Stability of training and generalization due to non-stationarity and out of distribution.

Computational cost (AutoRL) = cost (AutoML) x cost (RL)

Hyperparameter taxonomy: numerical, symbolic, neural. Stationary or nonstationary.

Learn selected parts of the system: learning HP, NAS, loss, rewards, agent and environment components.

Future: Co-dependencies. Nonstationary hyperparameters. Scaling up.

Demos: github.com/google/vizier/tree/main/demos github.com/google/brain_autorl/blob/main/evolving_rl/EvolvingRL_Demo.ipynb github.com/jcoreyes/evolvingrl