

CNNGen: Beyond Accuracy - A Generator and Benchmark for Energy-Aware Neural Architecture Search

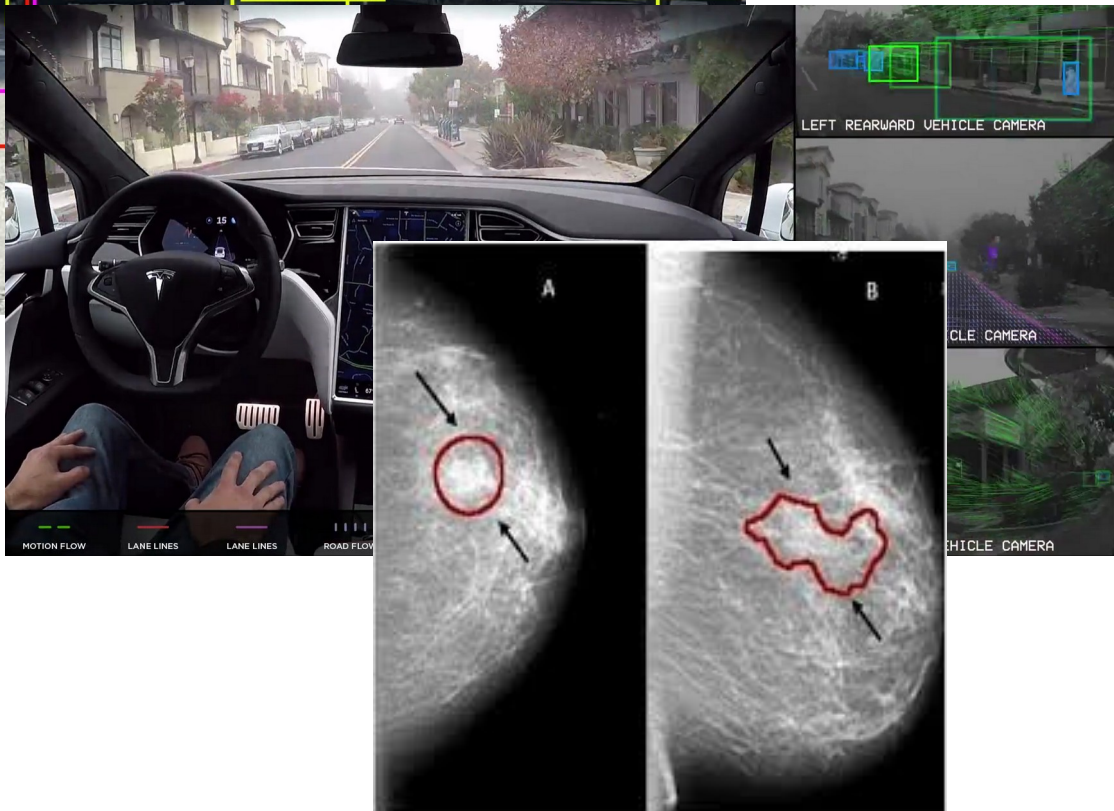
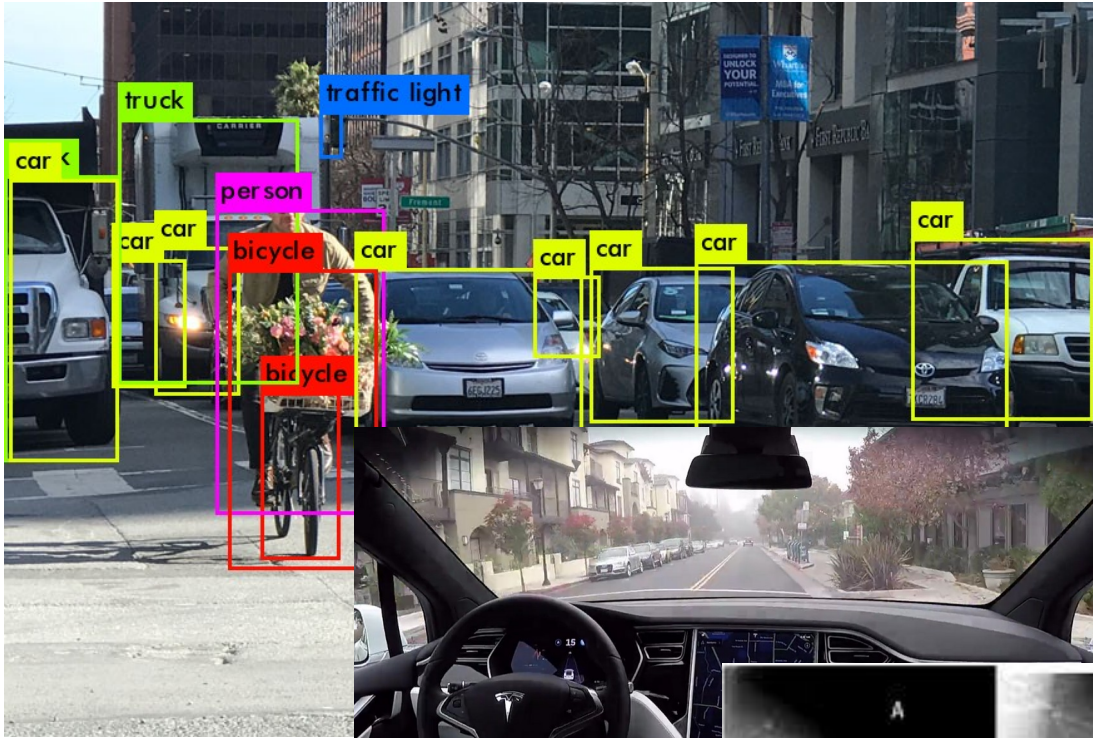
Antoine Gratia

DiverSE Coffee

Myself

- Antoine GRATIA
- PhD student
- University of Namur Belgium
- Supervisors
 - Gilles Perrouin (Unamur)
 - Pierre-Yves Schobbens (Unamur)





Context

- AI Race for Performance

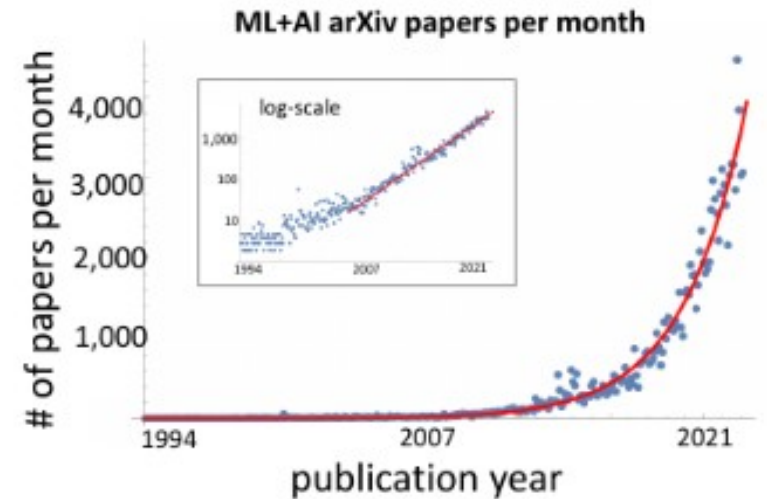
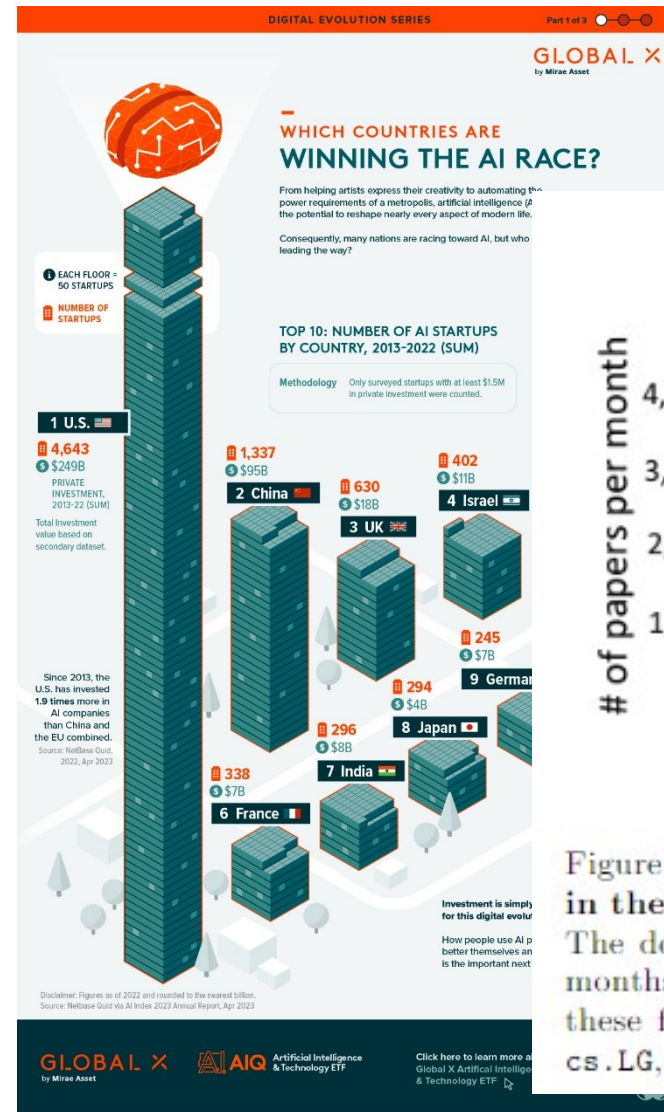


Figure 1. Number of papers published per months in the arXiv categories of AI grow exponentially. The doubling rate of papers per months is roughly 23 months, which might lead to problems for publishing in these fields, at some point. The categories are cs.AI, cs.LG, cs.NE, and stat.ML.

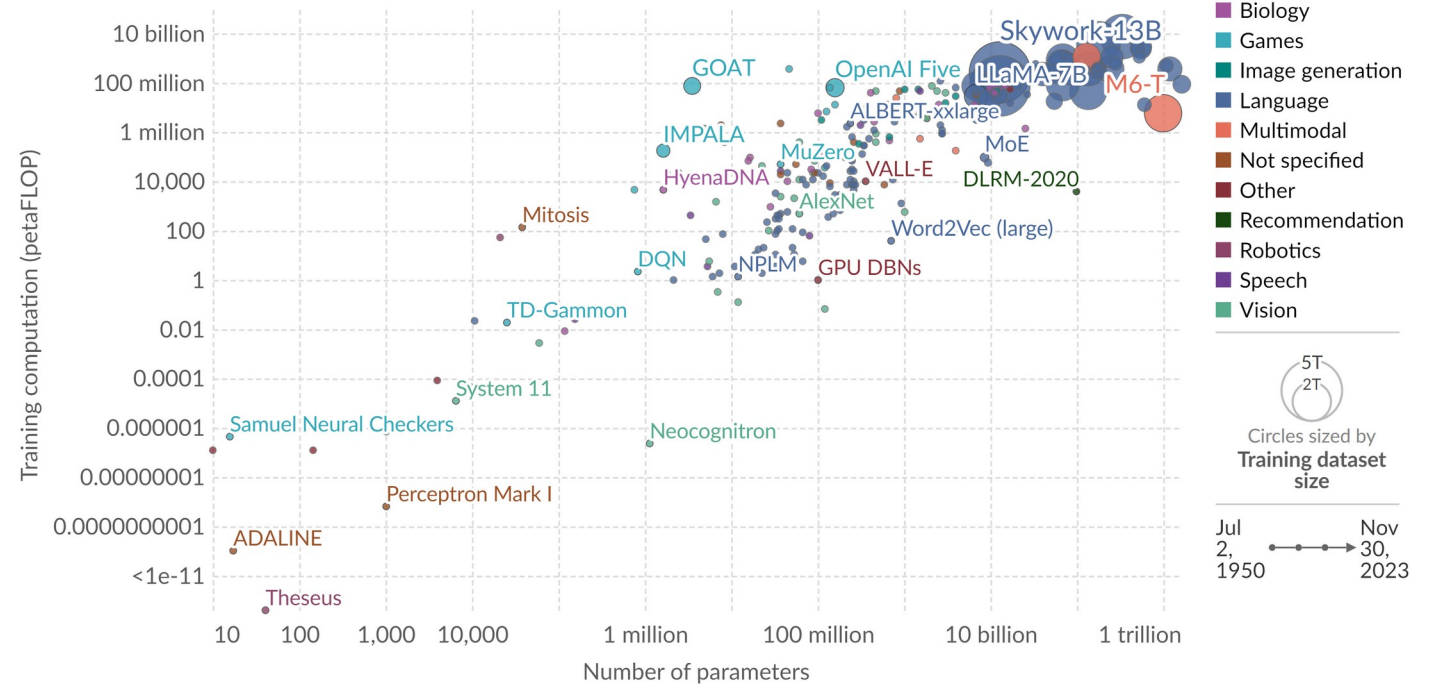
Energy

- LLaMA-65B cost 173 tCO2eq
- "AI me to the moon... Carbon footprint for training GTP-3 same as driving to our natural satellite and back"

Training computation vs. parameters in notable AI systems, by domain

Our World
in Data

Computation is measured in total petaFLOP, which is 10^{15} floating-point operations¹ estimated from AI literature, albeit with some uncertainty. Parameters are variables in an AI system whose values are adjusted during training to establish how input data gets transformed into the desired output.



Data source: Epoch (2023)

OurWorldInData.org/artificial-intelligence | CC BY

Note: Parameters are estimated based on published results in the AI literature and come with some uncertainty. The authors expect the estimates to be correct within a factor of 10.

1. **Floating-point operation:** A floating-point operation (FLOP) is a type of computer operation. One FLOP is equivalent to one addition, subtraction, multiplication, or division of two decimal numbers.

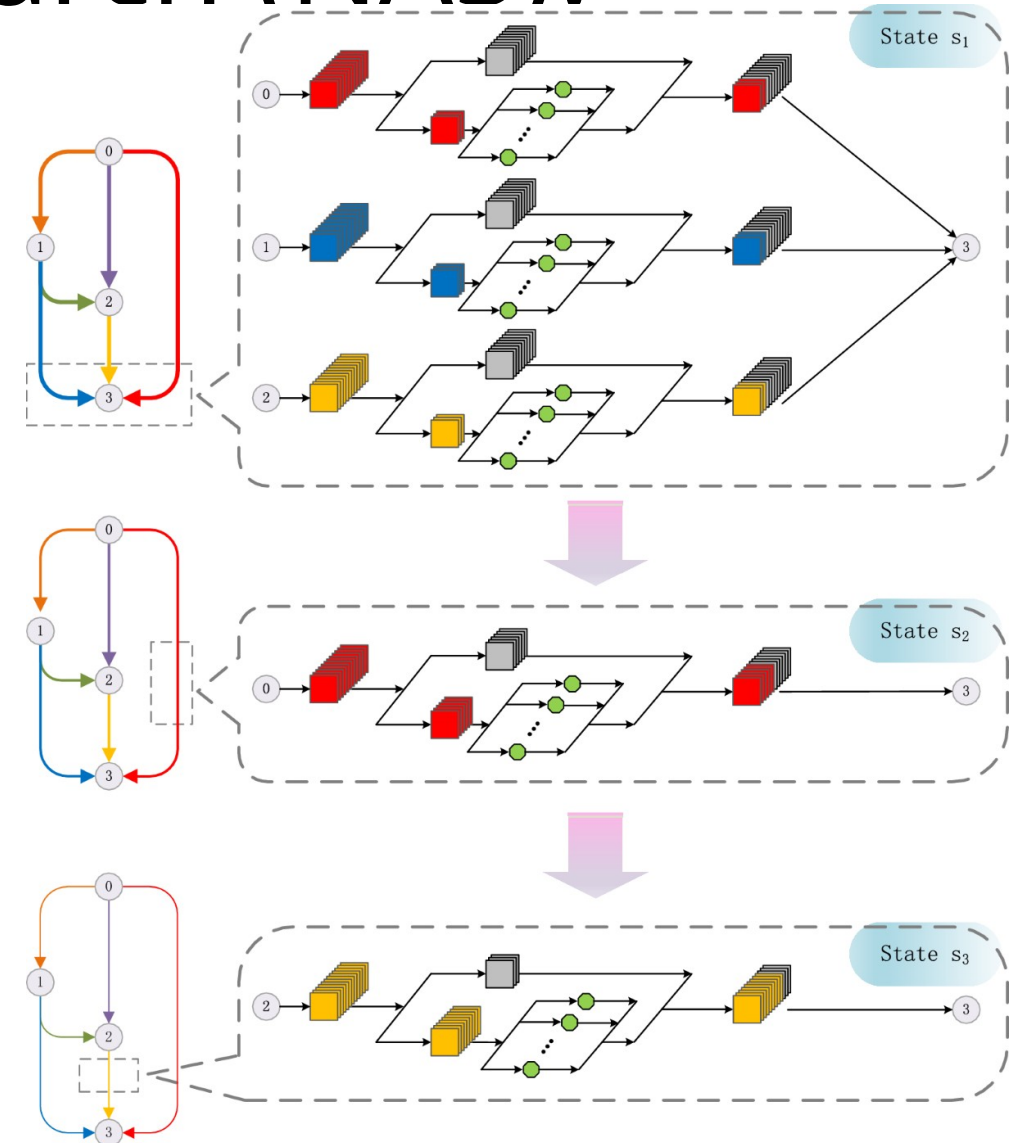
Objectives

Select architectures balancing performance & energy cost

- Model the search space
- Explore it (sampling)
- Energy cost and Performance Optimisation

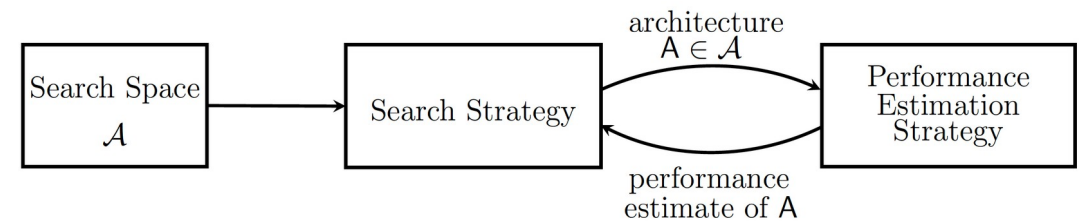
Neural Architecture Search (NAS)?

- Design neural network architectures automatically
- Find the optimal for a given task
- Exploring a predefined space (evolutionary algorithms, gradient-based optimisation,...)



NAS - General Workflow

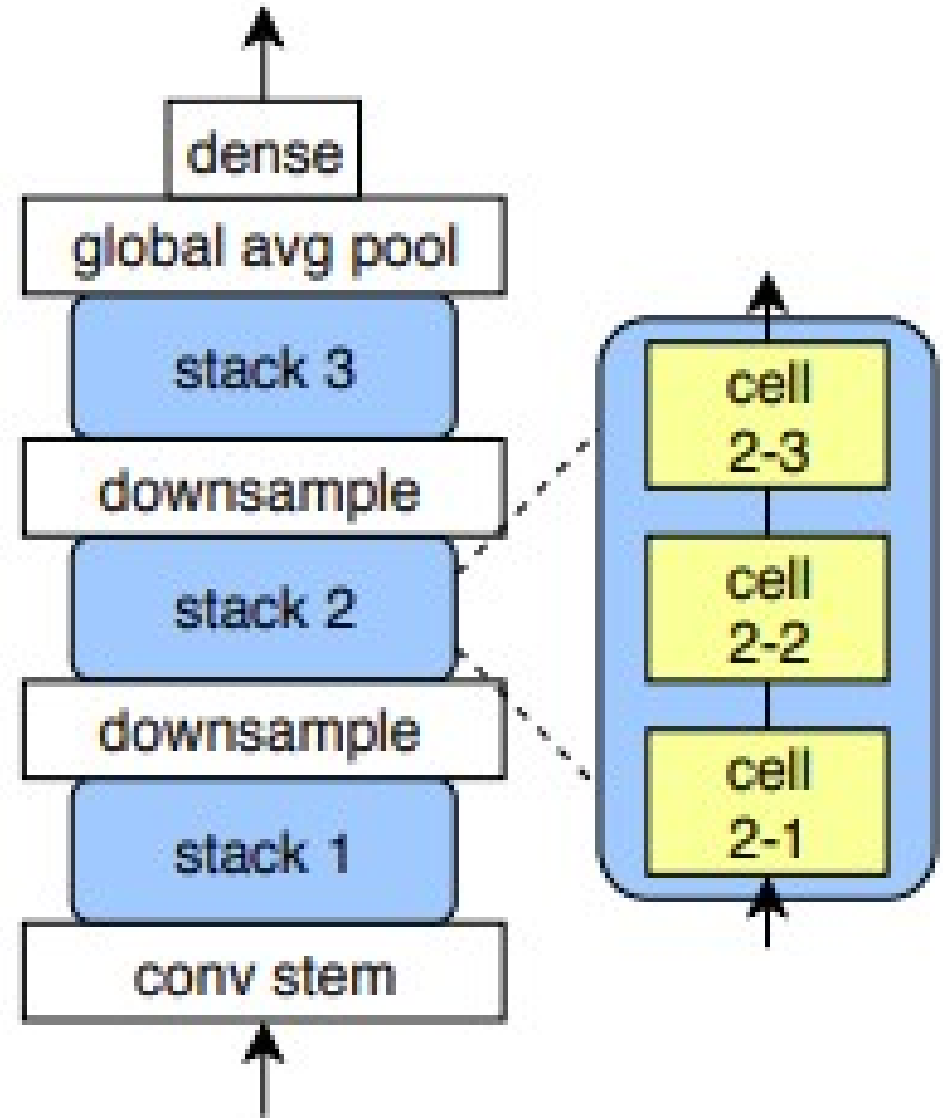
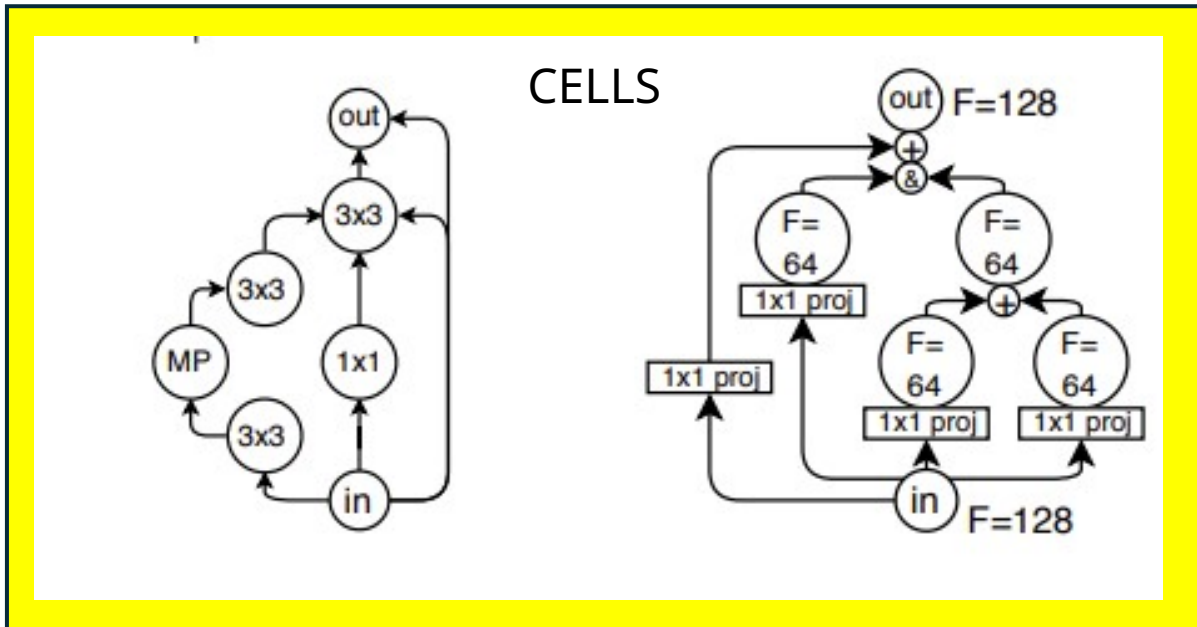
- Explores a search space
- Evaluates each architecture (train, prediction, ...)
- Estimates performance metric (accuracy, efficiency,...)



Goal: Find the optimal architecture for a given task

NAS – Example : NASBench-101

- Espace predifined
- 423K unique architectures
- Max 108 epochs



State Of The Art

Auto-Keras: An Efficient Neural Architecture Search System (Haifeng Jin - 06/2018)

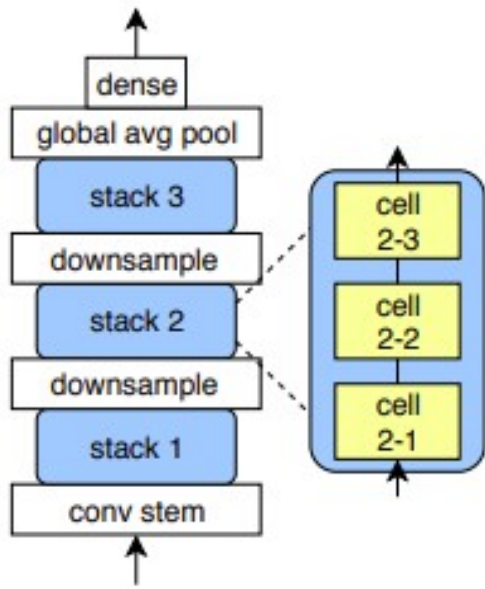
NATS-Bench: Benchmarking NAS Algorithms for Architecture Topology and Size
(Xuanyi Dong – 01/2021)

NAS-Bench-101: Towards Reproducible Neural Architecture Search (Chris Ying – 02/2019)

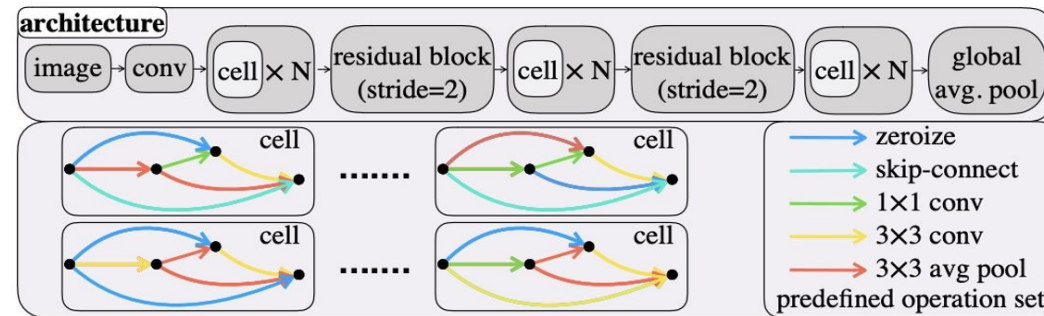
Construction of Hierarchical Neural Architecture Search Spaces based on Context-free Grammars (Simon Schrodi - 02/2024)

EA-HAS-Bench: Energy-Aware Hyperparameter and Architecture Search Benchmark (Shuguang Dou - 02/2023)

NAS - Problems



Pre-selection
(architecture)

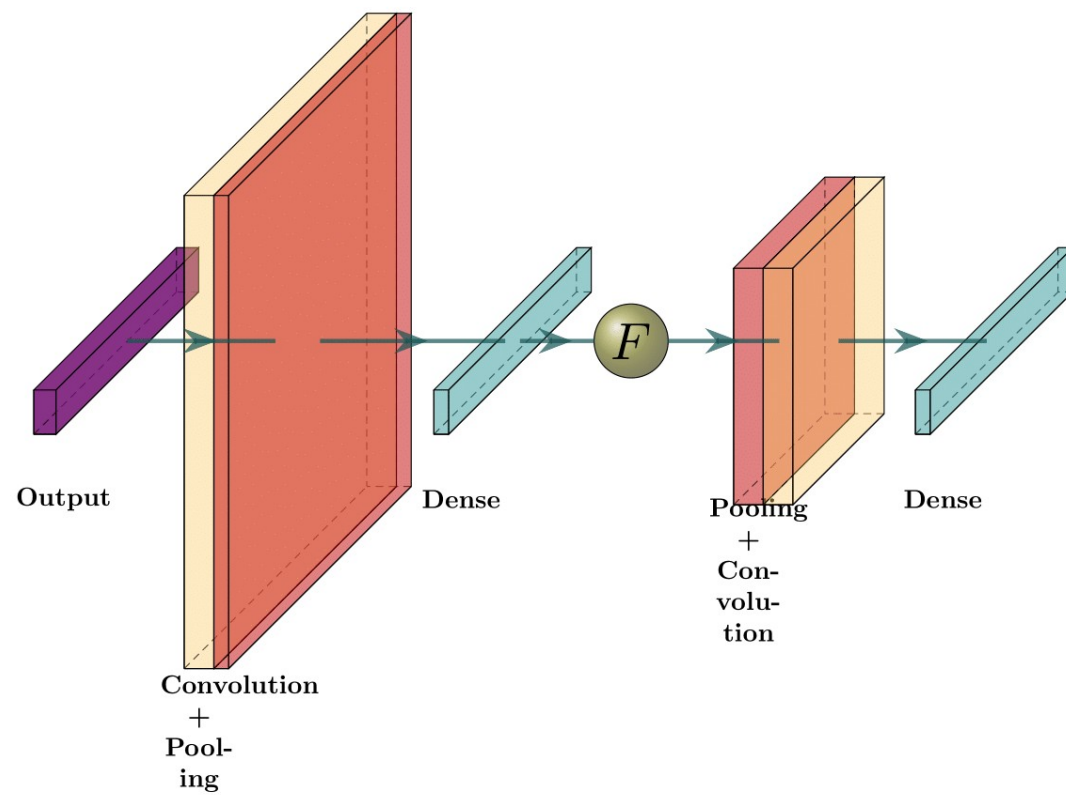
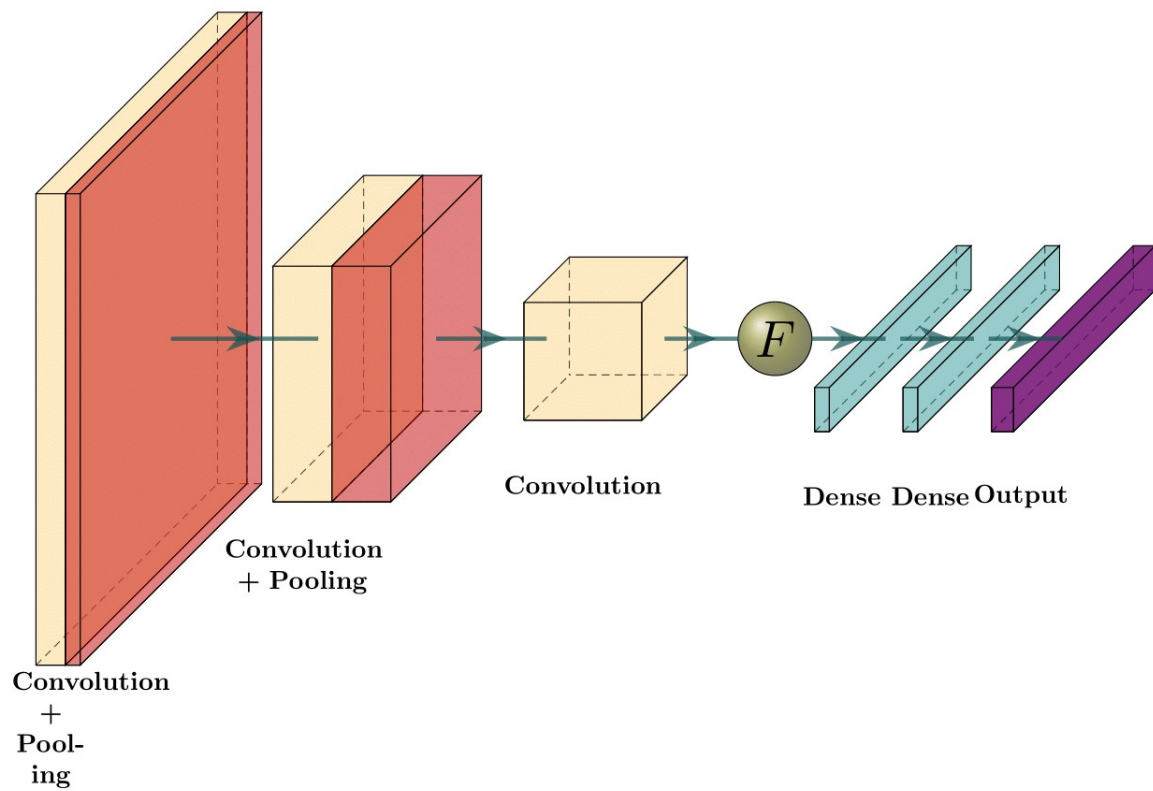


Variability (limited)

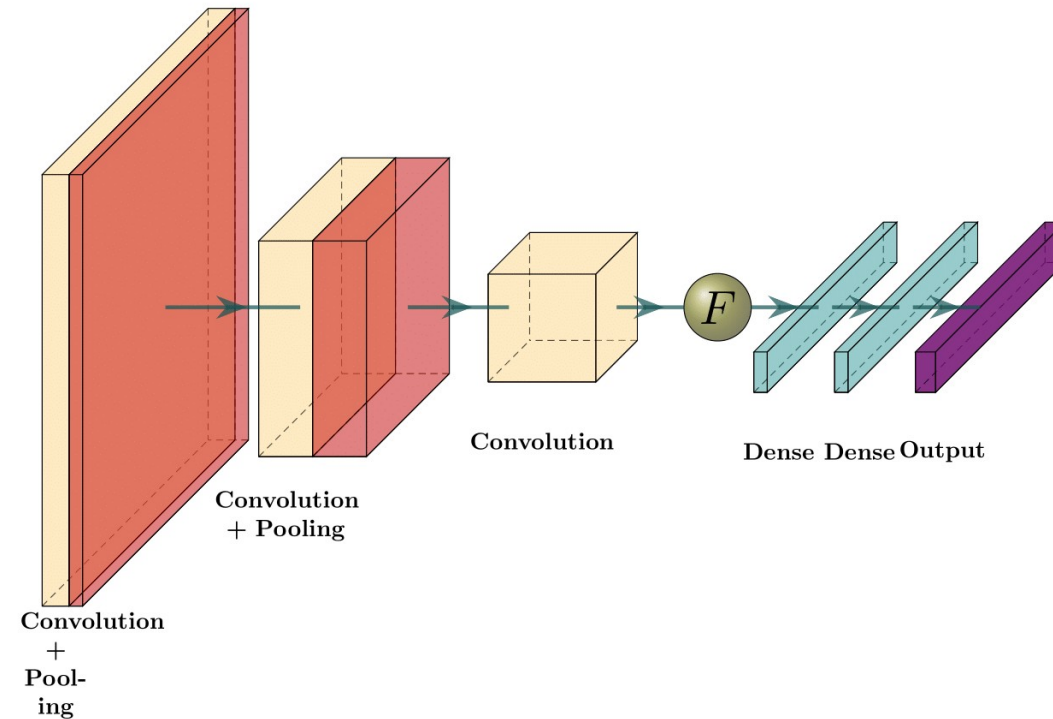
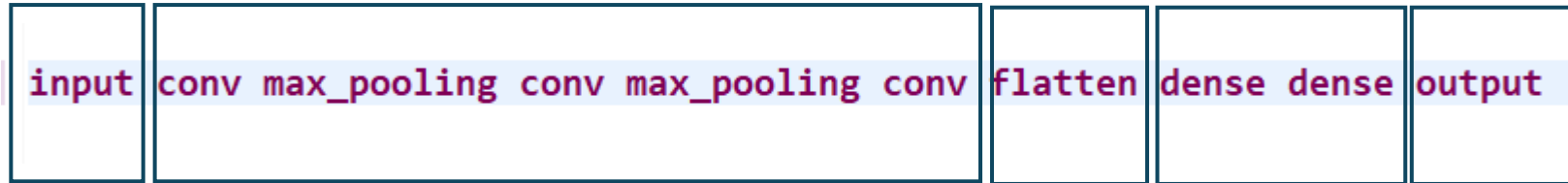


Need for
(re)training
(or pre-trained)





Example



Why?

```
input conv max_pooling conv max_pooling conv flatten dense dense output
```

Domain Specific Language

- CNN architectures' grammar
- Constraints support
- Validity check

Xtext
Xtend

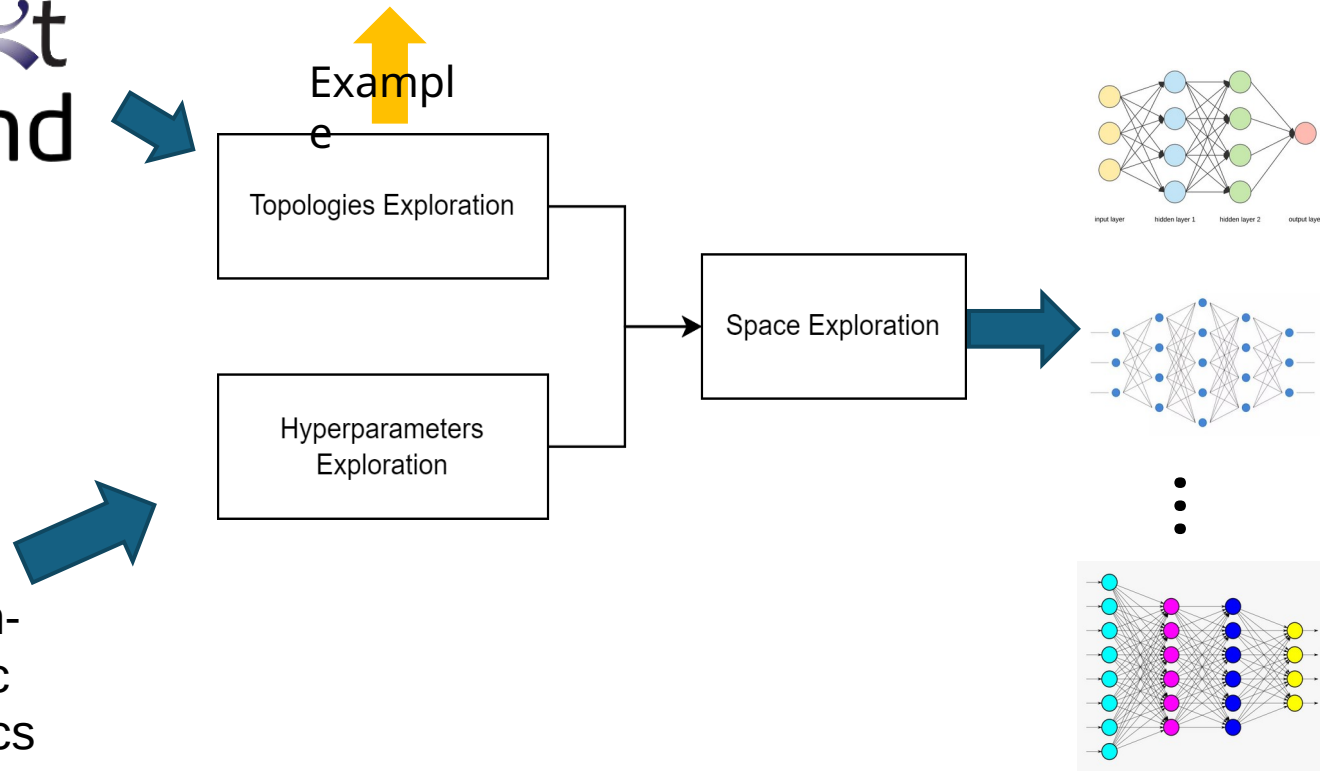
How did we find the grammar rules ?

- Based on domain knowledge
- Based on SOTA

CNNGEN

input conv max_pooling conv max_pooling conv flatten dense dense output

Xtext
Xtend



Selecting Candidate Architectures

Search strategies

- Random search
- Evolutionary algorithms

Select model to train

- Performance/energy prediction
- Deep Metric Learning
- Embedding based on grammar
- Distance (Levenshtein)
- Graph

Selecting Candidate Architectures

Search strategies

- Random search
- Evolutionary algorithms (Future work)

Select model to train

- **Performance/energy prediction**
- ~~Deep Metric Learning~~
- ~~Embedding based on grammar~~
- ~~Distance (Levenshtein)~~
- ~~Graph~~

State of the Art Performance Prediction NAS Overview

- **N-shot** evaluation method: The number of trained architectures \geq the number of searched architectures. All searched architectures are trained, and the focus is on accelerating their training to reduce runtime compared to Traditional Evaluation Method (TEM).
 - Downscaled dataset methods \Rightarrow reduce the dataset scale
 - Downscaled model methods \Rightarrow focus on reducing the model size
 - Network morphism and learning curve extrapolation \Rightarrow reduce the number of epochs

State of the Art Performance Prediction NAS Overview

- N-shot evaluation method
- **Few-shot** evaluation method: The number of trained architectures is less than the number of searched architectures, but greater than one. The runtime is naturally less than TEM due to fewer trained architectures.
 - Performance predictor
 - Population memory

State of the Art Performance Prediction NAS Overview

- N-shot evaluation method
- Few-shot evaluation method
- **One-shot** evaluation method: Only one architecture is trained, resulting in a runtime lower than TEM. The focus is on training a single architecture.
 - Path-based methods
 - Gradient-based methods

State of the Art Performance Prediction NAS Overview

- N-shot evaluation method
- Few-shot evaluation method
- One-shot evaluation method
- **Zero-shot** evaluation method: No architectures are trained, resulting in an extremely low cost. This method requires no training.
 - Parameter-level methods
 - Architecture-level methods

State of the Art Performance Prediction NAS Overview

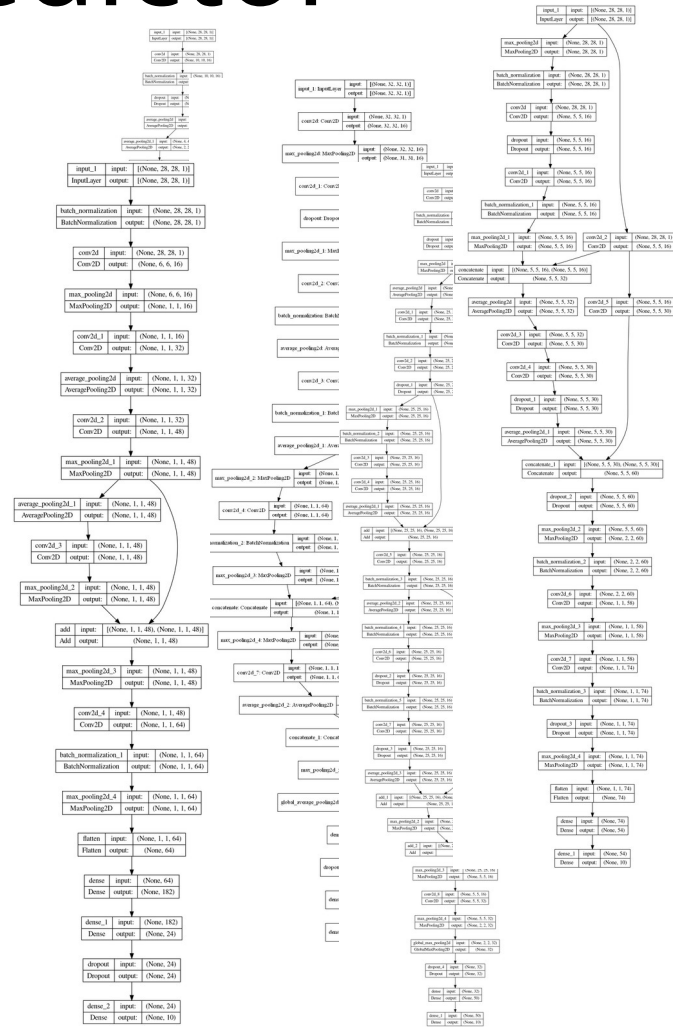
- N-shot evaluation method:
- **Few-shot evaluation method**
 - **Performance predictor**
 - Population memory
- One-shot evaluation method
- Zero-shot evaluation method

Predictor ¼ : Grammar-Based Predictor

input conv max_pooling conv max_pooling conv flatten dense dense output



Predictor 2/4 : TensorBoard Visualization P redictor



Predictor $\frac{3}{4}$: Code-Based Predictor

```
import numpy as np
import sys
from tensorflow.keras.callbacks import Callback, EarlyStopping, TensorBoard
from tensorflow.keras.models import Sequential, Model, load_model
from tensorflow.keras.layers import Input, Add, Dense, Activation, ZeroPadding2D, BatchNormalization, Flatten, Conv2D, AveragePooling2D,
from tensorflow.keras.initializers import glorot_uniform
from tensorflow.keras.utils import plot_model
from codecarbon import OfflineEmissionsTracker
import tensorflow as tf
import tensorflow as tf
import sys
import traceback
train_result_loss = ""
import csv
train_result_acc = ""
from time import time

# Load dataset
nb_layers = "not build"
(train_x, train_history) = None
epochs = 0

# normalize to 0
train_x = train_x
test_x = test_x
# start Emission Tracker
tracker = OfflineEmissionsTracker(country_iso_code="BEL", log_level="error", output_file="architecture_9_emissions.csv", output_dir="..")

val_x = train_x
val_y = train_y

try:
    def getModel():
        X_input = X - In try:
            X = Conv2D(16, k
            X = BatchNormali
            X = MaxPooling2D
            X = BatchNormali
            X = Conv2D(32, k
            X = AveragePooli
            X = Conv2D(48, k
            X = BatchNormali
            X = AveragePooli
            X = Conv2D(64, k
            X = BatchNormali
            X = MaxPooling2D
            X = Flatten()(X)
            X = Dense(112, a
            X = Dropout(0.88
            X = Dense(62, ac
            X = Dropout(0.18
            X = Dense(190, a
            model = Model(in
            return model
        model = getModel()

    def getModel():
        X_input = X = Input([32, 32, 3])
        X = Conv2D(16, kernel_size=3, strides=2, activation='selu', padding='valid')(X)
        X = BatchNormalization(epsilon=1e-5, axis=3)(X)
        X = MaxPooling2D(pool_size=2, strides=2, padding='same')(X)
        X1 = X
        X = AveragePooling2D(pool_size=3, strides=1, padding='same')(X)
        X = BatchNormalization(epsilon=1e-5, axis=3)(X)
        X = Conv2D(16, kernel_size=2, strides=1, activation='selu', padding='same')(X)
        X = Dropout(0.01)(X)
        X = Conv2D(16, kernel_size=2, strides=1, activation='selu', padding='same')(X)
        X = MaxPooling2D(pool_size=2, strides=1, padding='same')(X)
        X = Concatenate()([X, X1])
        X = MaxPooling2D(pool_size=1, strides=1, padding='valid')(X)
        X = BatchNormalization(epsilon=0.001, axis=3)(X)
        X = Conv2D(42, kernel_size=1, strides=1, activation='selu', padding='same')(X)
        X = MaxPooling2D(pool_size=4, strides=1, padding='valid')(X)
        X1 = X
        X = MaxPooling2D(pool_size=1, strides=1, padding='same')(X)
        X = Conv2D(42, kernel_size=2, strides=1, activation='selu', padding='same')(X)
        X = AveragePooling2D(pool_size=2, strides=1, padding='same')(X)
        X = Concatenate()([X, X1])
        X = AveragePooling2D(pool_size=1, strides=1, padding='valid')(X)
        X = BatchNormalization(epsilon=0.001, axis=3)(X)
        X = Conv2D(92, kernel_size=1, strides=1, activation='selu', padding='valid')(X)
        X = AveragePooling2D(pool_size=1, strides=1, padding='same')(X)
        X = GlobalAveragePooling2D()(X)
        X = Dropout(0.20)(X)
        X = Dense(114, activation='selu')(X)
        X = Dense(100, activation='softmax')(X)
        model = Model(inputs=X_input, outputs=X)
        return model
```



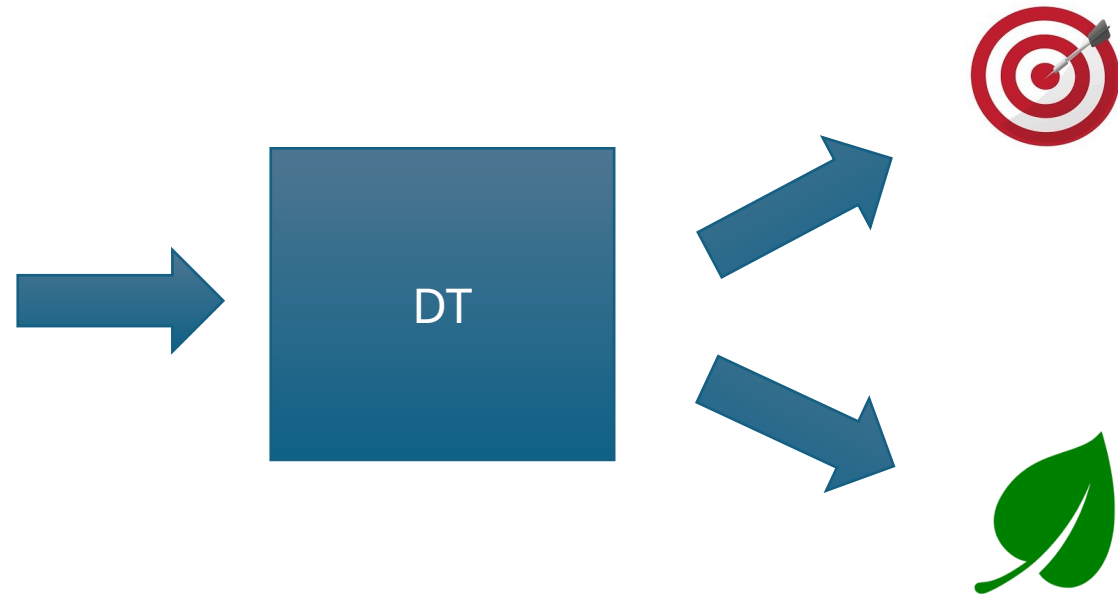
Transformers



Predictor 4/4 : Decision Tree-Based Predictor

Manually selected features:

- 0 Layers
- 0 Epochs
- 0 FLOP
- 0 Parameters



Experimentation

- Generate 1300 architectures with CNNGen
- Train Dataset (CIFAR-10, CIFAR-100, f-MNIST)
- Architecture Dataset to train predictors (energy & performance)

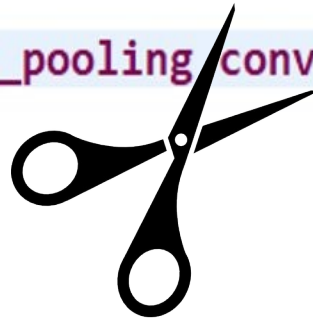
Results

- Decision Tree-Based Predictor => best result for energy prediction
- Code-Based Predictor => best result for performance prediction

Future works - Neuro-evolution

- Automatic optimisation => genetic algorithms
 - JMetal
 - Genetic improvement on grammar instance

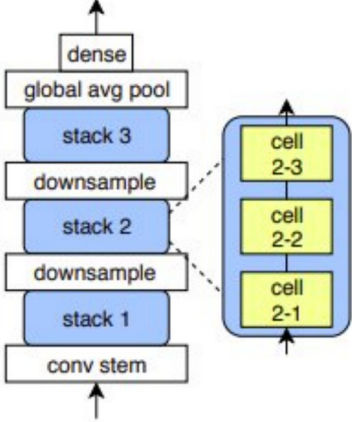
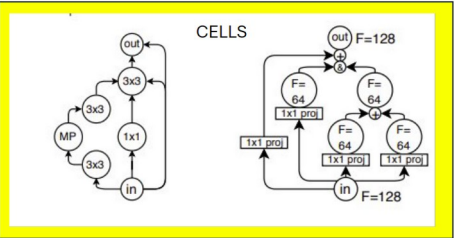
```
input conv max_pooling conv max_pooling conv flatten dense dense output
```



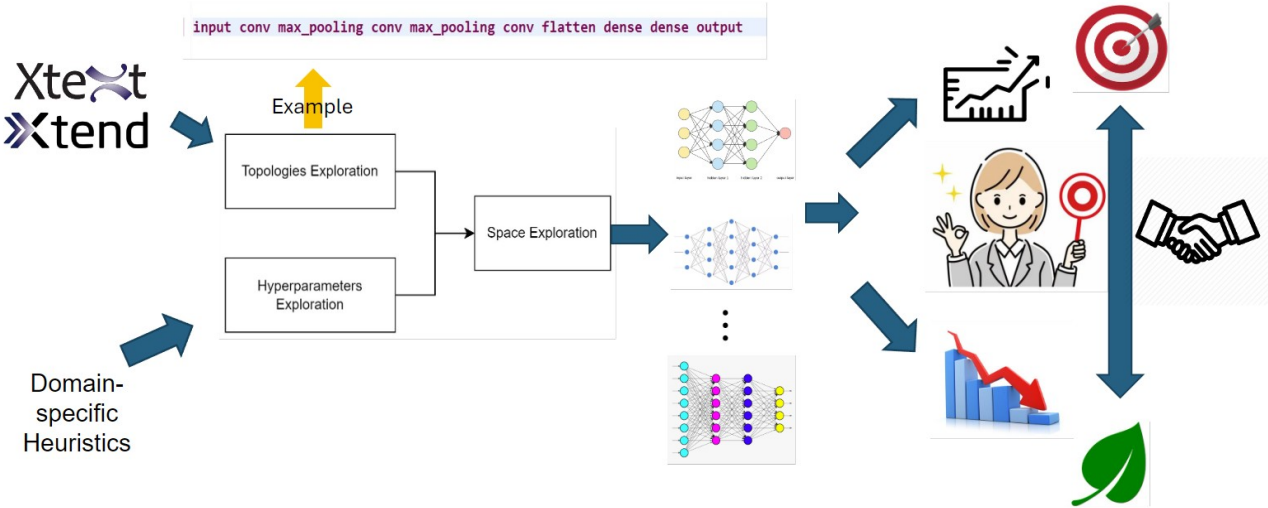
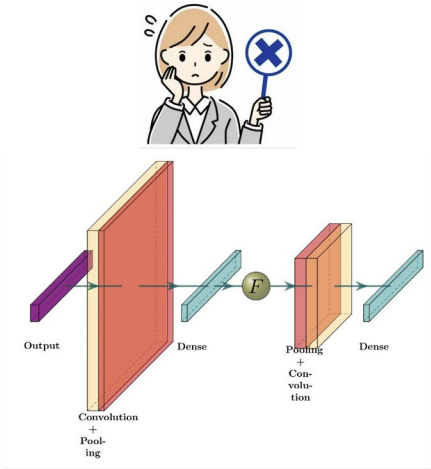
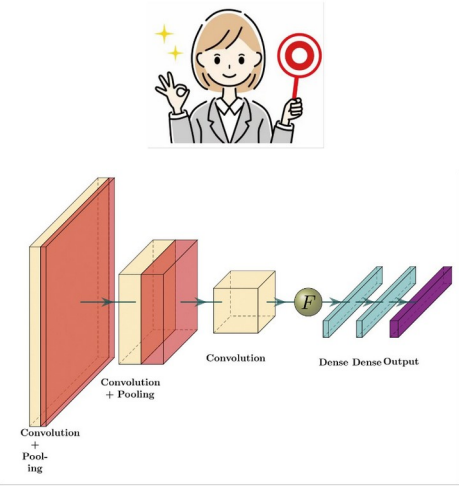
Sum up

NAS - Example

- Espace predefined
- 423K unique architectures
- Max 108 epochs



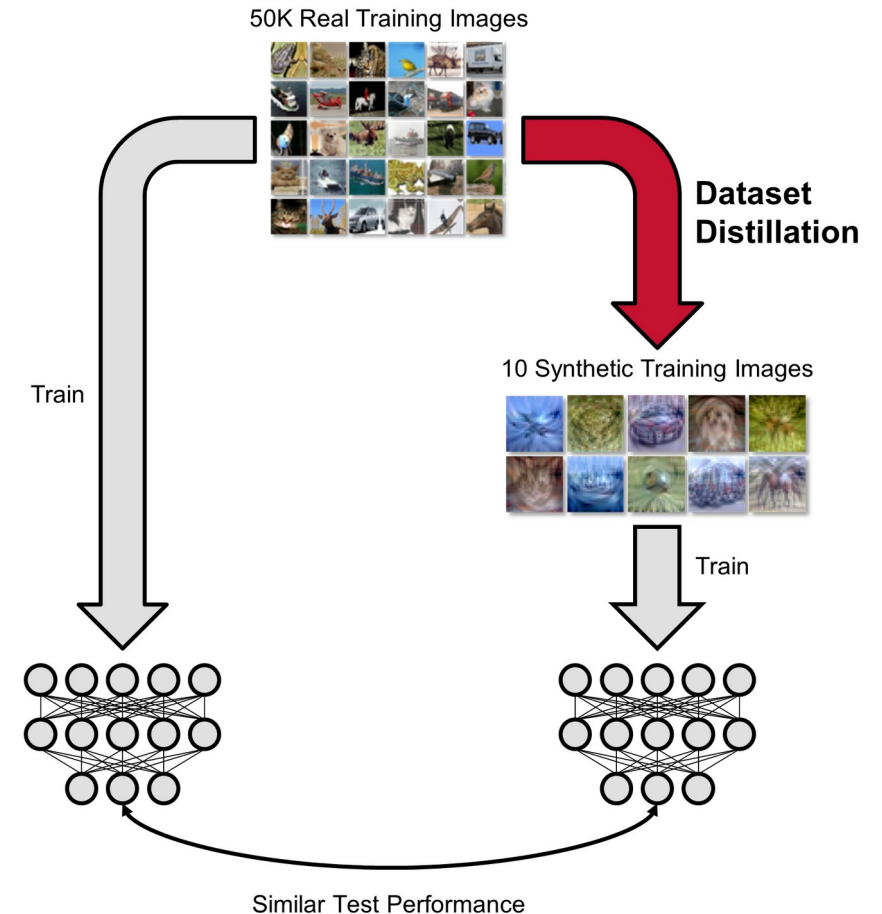
CNNGEN



BackUP

Future works - Dataset Distillation

- Create SynD
 - Same performance
 - Less cost
- Benchmark divers DD techniques (Energy cost)
- Train selected CNNGen architectures



The rise of artificial intelligence over the last 8 decades: As training computation has increased, AI systems have become more powerful

The color indicates the domain of the AI system: ● Vision ● Games ● Drawing ● Language ● Other

Shown on the vertical axis is the **training computation** that was used to train the AI systems.



The data on training computation is taken from Sevilla et al. (2022) – Parameter, Compute, and Data Trends in Machine Learning. It is estimated by the authors and comes with some uncertainty. The authors expect the estimates to be correct within a factor of two. OurWorldInData.org – Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the authors Charlie Giattino, Edouard Mathieu, and Max Roser

State of the Art Performance Prediction NAS Overview

- N-shot evaluation method: The number of trained architectures is greater than or equal to the number of searched architectures. All searched architectures are trained, and the focus is on accelerating their training to reduce runtime compared to Traditional Evaluation Method (TEM).
 - Downscaled dataset methods => reduce the dataset scale
 - Downscaled model methods => focus on reducing the model size
 - Network morphism and learning curve extrapolation => reduce the number of epochs
- Few-shot evaluation method: The number of trained architectures is less than the number of searched architectures, but greater than one. The runtime is naturally less than TEM due to fewer trained architectures.
 - Performance predictor
 - Population memory
- One-shot evaluation method: Only one architecture is trained, resulting in a runtime lower than TEM. The focus is on training a single architecture.
 - Path-based methods
 - Gradient-based methods
- Zero-shot evaluation method: No architectures are trained ($A_t = \emptyset$), resulting in an extremely low cost. This method requires no training.
 - Parameter-level methods
 - Architecture-level methods

N-shot evaluation method

- Downscaled dataset methods
 - Reduced Training Time & Accelerated Exploration
 - Using a smaller dataset, the training process becomes quicker and more efficient (beneficial large and complex datasets)
 - Allow to iterate more rapidly in architecture search and experimentation.
- Downscaled model methods
 - Reduced Training Time & Accelerated Exploration
 - Reducing the size of the models the training process becomes faster and more efficient (beneficial large and complex models)
 - Allow to iterate more rapidly in architecture search and experimentation.
- Network morphism
- Learning curve extrapolation

N-shot evaluation method

- Downscaled dataset methods
- Downscaled model methods
- Network morphism
 - Accelerated Training Process
 - Child network can inherit the weights and function of a fully-trained parent network
=> child network does not need to be trained from scratch and requires fewer epochs to converge
- Learning curve extrapolation
 - Reduced training time
 - Predict the final performance of architectures trained for only a few epochs
 - Incorporation of Partial Learning Curves:
 - Partial learning curves captures valuable information about the architecture's learning progress => accurate predictions of the architecture performance, even with limited training epochs.

Few-shot evaluation method

- Performance predictor
 - Reduced runtime
 - Estimate the performance of architectures in the search space without actually training and evaluating
 - Resource efficiency
 - Build a small architecture set (A_t) that represents the search space adequately => reduces the number of architectures that need to be trained and evaluated
 - Flexibility and adaptability
 - Updated during the architecture search process, allowing them to incorporate new samples and improve their prediction performance.
 - Add new architectures to the training set and retrain the predictor, making them more flexible and adaptable to changing search spaces.
- Population Memory
 - Accelerated Fitness Evaluation
 - Avoiding the repeated evaluation of the same architecture by reusing architectural information that has appeared in previous populations(as a cache system)

One-shot evaluation method

- Path-based methods
 - Accelerated performance evaluation
 - Instead of training and evaluating each subnet individually, the weights of the subnets are directly extracted from the supernet, and their performance is inferred on a validation dataset
 - Weight decoupling
 - To decouple the weights from the architecture and improve the performance evaluation sampling strategies are employed during the supernet training
 - Optimize the weights of different subnets in a more balanced way, effectively alleviating the weight coupling problem and providing better performance estimation for the subnets.
- Gradient-based methods
 - Decoupling of training and architecture search
 - Separate the training of the supernet and the search for the architectures => allows for more efficient optimization by jointly optimizing the supernet weights and architecture parameters
 - Continuous relaxation of the search space
 - Relax the discrete search space to be continuous => softmax function allows for the continuous representation of the architecture choices
 - Relaxation enables to efficiently search for optimal architectures
 - Challenges
 - Huge memory consumption
 - Poor generalization
 - Performance collapse

Zero-shot evaluation method

- Parameter-level methods
 - Reduced computational cost
 - Require a minibatch of data and a forward/backward propagation pass to calculate saliency indicators for certain parameters.
 - Don't involve training the architecture, the computational requirements are significantly reduced.
 - Utilization of network pruning techniques:
 - Leverage saliency indicators (used in network pruning literature) measure the importance of parameters and help identify unimportant ones
 - Can evaluate the saliency of parameters and score the entire architecture based on the aggregated saliencies.
- Architecture-level methods
 - Time-saving performance evaluation
 - Evaluate the performance of architectures by measuring properties related to architecture performance without training
 - Theoretical basis and indicators
 - Indicators as activation overlap, learnability, gradient behavior, expressivity, synaptic diversity, ... to judge the properties of architectures and rank them
- Inaccurate => performance fluctuated dramatically among different tasks.