# Thank you, **tinyML Strategic Partners**, for committing to take tinyML to the next Level, together

# Executive Strategic Partners

Qualcomm
AI research

# Advancing AI research to make efficient AI ubiquitous

**Power efficiency**

Model design, compression, quantization, algorithms, efficient hardware, software tool

**Personalization**

Continuous learning, contextual, always-on, privacy-preserved, distributed learning

**Efficient learning**

Robust learning through minimal data, unsupervised learning, on-device learning

# A platform to scale AI across the industry

**Perception**
Object detection, speech recognition, contextual fusion

**Reasoning**
Scene understanding, language understanding, behavior prediction

**Action**
Reinforcement learning for decision making

IoT/IIoT

Edge cloud

Automotive

Cloud

Mobile

# Platinum Strategic Partners

Renesas is enabling the next generation of AI-powered solutions that will revolutionize every industry sector.

AGRICULTURE

BUILDING AUTOMATION

CITY INFRASTRUCTURE

TRANSPORT & LOGISTIC

HEALTH CARE

HOME AUTOMATION

INDUSTRIAL AUTOMATION

CONSUMERS

renesas.com

RENESAS
BIG IDEAS FOR EVERY SPACE

**DEPLOY VISION AI AT THE EDGE AT SCALE**

SONY

Gold Strategic Partners

ARDUINO® PRO™

Easily deploy your
tinyML solutions with
Arduino Pro

arduino.cc/pro

Made In Italy

Powering tinyML Innovation

# Arm AI Virtual Tech Talks

The latest in AI trends, technologies & best practices from Arm and our Ecosystem Partners.

Demos, code examples, workshops, panel sessions and much more!

Fortnightly Tuesday @ 4pm GMT/8am PT

**Find out more:**
**www.arm.com/techtalks**

# NEUROMORPHIC INTELLIGENCE FOR THE SENSOR-EDGE

www.innatera.com

Microsoft

# The Right Edge AI Tools Can Make or Break Your Next Smart IoT Product

## Analytics Toolkit Suite

**SensiML**

AutoML

Data Collection

Test & Validation

Data Labeling

Code Generation

Model Building

Team Collaboration

Version Control and Model Management

sensiml.com/tinyML

STMicroelectronics provides extensive solutions to make tiny Machine Learning easy

life.augmented

www.st.com/ai

# Silver Strategic Partners

# Join Growing tinyML Communities:



13.9k members in
47 Groups in 39 Countries

**tinyML - Enabling ultra-low Power ML at the Edge**
https://www.meetup.com/tinyML-Enabling-ultra-low-Power-ML-at-the-Edge/



4k members
&
11.6k followers

**The tinyML Community**
https://www.linkedin.com/groups/13694488/

**Subscribe to
tinyML YouTube Channel
for updates and notifications
*(including this video)***

www.youtube.com/tinyML

*EMEA 2023*

https://www.tinyml.org/event/emea-2023

More sponsorships are available: sponsorships@tinyML.org

# Reminders

Slides & Videos will be posted tomorrow

tinyml.org/forums     youtube.com/tinyml

Please use the Q&A window for your questions

# Swarnava Dey



Swarnava Dey is a Senior Scientist at TCS Research working on embedded vision systems. He is an M.Tech from IIT, Kharagpur, and currently pursuing PhD there in robustness, verifiability and explainability of Embedded Deep Neural Networks and Neuro Symbolic AI. He has 30+ granted patents, 25+ research papers, and is an author of Towards Data Science: https://medium.com/@qswadey. His publication details can be found at his Google Scholar page:https://scholar.google.co.in/citations?hl=en&user=aFplwjEAAAAJ

# Neural Architecture Search (~AutoML): Goals

- <u>NAS research</u> - Automatically generate better architectures than handcrafted models; benchmark accuracy on NAS Bench dataset & ImageNet

- <mark><u>NAS for TinyML</u></mark> - Automatically customize & optimize DNNs for multiple constraints - [accuracy, model size, SRAM usage (runtime memory), #MACs, latency, energy usage…]

tcs Research

Inventing for impact

# Neural Architecture Search (~AutoML): Goals

- <u>NAS research</u> - Automatically generate better architectures than handcrafted models; benchmark accuracy on NAS Bench dataset & ImageNet

- <mark>NAS for TinyML</mark> - Automatically customize & optimize DNNs for multiple constraints - [accuracy, model size, SRAM usage (runtime memory), #MACs, latency, energy usage…]

💡 NAS SOTA - Differential, one-shot NAS - DARTS & DARTS-based
- ✔ DARTS works out of the box - https://github.com/quark0/darts
- ✔ Default implementation single-objective - accuracy
- ✔ Difficult to get ready version that allows integration of my preferred objectives

tcs Research

Inventing for impact

# Neural Architecture Search (~AutoML): Goals

- <u>NAS research</u> - Automatically generate better architectures than handcrafted models; benchmark accuracy on NAS Bench dataset & ImageNet

- <mark>NAS for TinyML</mark> - Automatically customize & optimize DNNs for multiple constraints - [accuracy, model size, SRAM usage (runtime memory), #MACs, latency, energy usage…]

👂 NAS SOTA - Differential, one-shot NAS - DARTS & DARTS-based
  - ✔ DARTS works out of the box - https://github.com/quark0/darts
  - ✔ Default implementation single-objective - accuracy
  - ✔ Difficult to get ready version that allows integration of my preferred objectives

👂 Heuristic / Sample based NAS are obsolete
  - ✔ For searching tiny models the overhead is marginalized
  - ✔ Easy to integrate preferred objectives

tcs Research

# Neural Architecture Search (~AutoML): Goals

- <u>NAS research</u> - Automatically generate better architectures than handcrafted models; benchmark accuracy on NAS Bench dataset & ImageNet

- <mark>NAS for TinyML</mark> - Automatically customize & optimize DNNs for multiple constraints - [accuracy, model size, SRAM usage (runtime memory), #MACs, latency, energy usage…]

👂 NAS SOTA - Differential, one-shot NAS - DARTS & DARTS-based
- ✔ DARTS works out of the box - https://github.com/quark0/darts
- ✔ Default implementation single-objective - accuracy
- ✔ Difficult to get ready version that allows integration of my preferred objectives

👂 Heuristic / Sample based NAS are obsolete
- ✔ For searching tiny models the overhead is marginalized
- ✔ Easy to integrate preferred objectives

👂 NAS for TinyML
- ✔ Large many-layered networks, complex connections not required
- ✔ Accurate multi-objective conformance, platform API support - highly required
- ✔ MCUNet V1 & V2, µ-NAS, Micronets… have their merits & demerits. We may often need to tweak existing frameworks - <u>understanding NAS helps</u>

**tcs Research**

Inventing for impact

# Today's Agenda

- SOTA NAS research & TinyML

- A naive method to generate DNN

- Enhancing the method using Reinforcement Learning

- Other optimization techniques for sample-based NAS

- Gradient-based, one-shot NAS

- Take-home points

tcs Research

Inventing for impact

# NAS Research & TinyML: Mainstream NAS Goals

▶ Started with RL using RNN[1] , Q-Learning[2] to generate architecture

**[1]** **Neural Architecture Search with Reinforcement Learning,** Barret Zoph, Quoc V. Le, https://arxiv.org/abs/1611.01578

**[2]** **Designing Neural Network Architectures using Reinforcement Learning**, Bowen Baker *et al.*, https://arxiv.org/abs/1611.02167

tcs Research

*Inventing for impact*

# NAS Research & TinyML: Mainstream NAS Goals

▶ Started with RL using RNN[1] , Q-Learning[2] to generate architecture

▶ Next two big innovations were cellular search space[3] and parameter sharing[4]

[1] **Neural Architecture Search with Reinforcement Learning,** Barret Zoph, Quoc V. Le, https://arxiv.org/abs/1611.01578
[2] **Designing Neural Network Architectures using Reinforcement Learning**, Bowen Baker *et al.*, https://arxiv.org/abs/1611.02167
[3] **Learning Transferable Architectures for Scalable Image Recognition,** Barret Zoph *et al.*, https://arxiv.org/abs/1707.07012
[4] **Efficient Neural Architecture Search via Parameter Sharing**, Hieu Pham *et al.*, https://arxiv.org/pdf/1802.03268.pdf

tcs Research

Inventing for impact

# NAS Research & TinyML: Mainstream NAS Goals

▶ Started with RL using RNN[1] , Q-Learning[2] to generate architecture

▶ Next two big innovations were cellular search space[3] and parameter sharing[4]

▶ Meanwhile Evolutionary[5], Bayesian NAS[6] pushed the NAS optimization SOTA

[1] **Neural Architecture Search with Reinforcement Learning,** Barret Zoph, Quoc V. Le, https://arxiv.org/abs/1611.01578

[2] **Designing Neural Network Architectures using Reinforcement Learning**, Bowen Baker *et al.*, https://arxiv.org/abs/1611.02167

[3] **Learning Transferable Architectures for Scalable Image Recognition,** Barret Zoph *et al.*, https://arxiv.org/abs/1707.07012

[4] **Efficient Neural Architecture Search via Parameter Sharing**, Hieu Pham *et al.*, https://arxiv.org/pdf/1802.03268.pdf

[5] **Regularized Evolution for Image Classifier Architecture Search**,Esteban Real *et al.*

[6] **Parallelised Bayesian Optimisation via Thompson Sampling,** Kirthevasan Kandasamy *et al.*

tcs Research

Inventing for impact

# NAS Research & TinyML: Mainstream NAS Goals

▶ Started with RL using RNN[1] , Q-Learning[2] to generate architecture

▶ Next two big innovations were cellular search space[3] and parameter sharing[4]

▶ Meanwhile Evolutionary[5], Bayesian NAS[6] pushed the NAS optimization SOTA

▶ **Next came the game changer DARTS[7]. It used earlier innovations.**

**[1] Neural Architecture Search with Reinforcement Learning,** Barret Zoph, Quoc V. Le, https://arxiv.org/abs/1611.01578
**[2] Designing Neural Network Architectures using Reinforcement Learning**, Bowen Baker *et al.*, https://arxiv.org/abs/1611.02167
**[3] Learning Transferable Architectures for Scalable Image Recognition,** Barret Zoph *et al.*, https://arxiv.org/abs/1707.07012
**[4] Efficient Neural Architecture Search via Parameter Sharing**, Hieu Pham *et al.*, https://arxiv.org/pdf/1802.03268.pdf
**[5] Regularized Evolution for Image Classifier Architecture Search**,Esteban Real *et al.*
**[6] Parallelised Bayesian Optimisation via Thompson Sampling,** Kirthevasan Kandasamy *et al.*
**[7] DARTS: Differentiable Architecture Search**, Hanxiao Liu *et al.*, https://arxiv.org/abs/1806.09055

tcs Research

Inventing for impact

# NAS Research & TinyML: Mainstream NAS Goals

- Started with RL using RNN[1] , Q-Learning[2] to generate architecture
- Next two big innovations were cellular search space[3] and parameter sharing[4]
- Meanwhile Evolutionary[6], Bayesian NAS[7] pushed the NAS optimization SOTA
- Next came the game changer DARTS[5]. It used earlier innovations.
- All paved the way for searching larger, more accurate nets much faster
- I had to drop many deserving papers due to space constraints
- This goal is not important for TinyML applications.
- However, let's remember **2**-MetaQNN, **3**-NASNet, **4**-ENAS, **5**-DARTS, **6**-AE, **& 7**-BO

[1] **Neural Architecture Search with Reinforcement Learning,** Barret Zoph, Quoc V. Le, https://arxiv.org/abs/1611.01578
[2] **Designing Neural Network Architectures using Reinforcement Learning**, Bowen Baker *et al.*, https://arxiv.org/abs/1611.02167
[3] **Learning Transferable Architectures for Scalable Image Recognition,** Barret Zoph *et al.*, https://arxiv.org/abs/1707.07012
[4] **Efficient Neural Architecture Search via Parameter Sharing**, Hieu Pham *et al.*, https://arxiv.org/pdf/1802.03268.pdf
[5] **Regularized Evolution for Image Classifier Architecture Search**,Esteban Real *et al.*
[6] **Parallelised Bayesian Optimisation via Thompson Sampling,** Kirthevasan Kandasamy *et al.*
[7] **DARTS: Differentiable Architecture Search**, Hanxiao Liu *et al.*, https://arxiv.org/abs/1806.09055

tcs Research

Inventing for impact

# NAS Research & TinyML: TinyML Goals

▶ 1000-layer DNN for MCUs ? - 20-30 layer DNN fits ; Global search space good enough; complex, unexpected interconnect[8] OK ; Training smaller models marginalizes many problems



**[8] SwiftNet: Using Graph Propagation as Meta-knowledge to Search Highly Representative Neural Architectures**, Cheng *et al.*

tcs Research

Inventing for impact

# NAS Research & TinyML: TinyML Goals

▶ 1000-layer DNN for MCUs ? - 20-30 layer DNN fits ; Global search space good enough; complex, unexpected interconnect[8] OK ; Training smaller models marginalizes many problems

▶ For TinyML model size, peak RAM usage, #MAC are highly important

▶ Evaluation, measurement of KPIs, reward engineering are important

**[8] SwiftNet: Using Graph Propagation as Meta-knowledge to Search Highly Representative Neural Architectures,** Cheng *et al.*

tcs Research

Inventing for impact

# NAS Research & TinyML: TinyML Goals

▶ 1000-layer DNN for MCUs ? - 20-30 layer DNN fits ; Global search space good enough; complex, unexpected interconnect[8] OK ; Training smaller models marginalizes many problems

▶ For TinyML model size, peak RAM usage, #MAC are highly important

▶ Evaluation, measurement of KPIs, reward engineering are important

▶ MCUNet v1[9] & v2[10] are top class works - based on one shot NAS[11], Search Space reduction highest FLOPs within SRAM constraint - closely tied with TinyEngine IR for HW params

[8] **SwiftNet: Using Graph Propagation as Meta-knowledge to Search Highly Representative Neural Architectures,** Cheng *et al.*
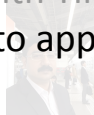
[9] **MCUNet: Tiny Deep Learning on IoT Devices**, Ji Lin *et al.*, https://arxiv.org/pdf/2007.10319

[10] **MCUNetV2: Memory-Efficient Patch-based Inference for Tiny Deep Learning**, Ji Lin *et al.*, https://arxiv.org/abs/2110.15352

[11] **Understanding and Simplifying One-Shot Architecture Search**, Bender *et al.*, https://proceedings.mlr.press/v80/bender_18a

tcs Research

Inventing for impact

# NAS Research & TinyML: TinyML Goals

▶ 1000-layer DNN for MCUs ? - 20-30 layer DNN fits ; Global search space good enough; complex, unexpected interconnect[8] OK ; Training smaller models marginalizes many problems

▶ For TinyML model size, peak RAM usage, #MAC are highly important

▶ Evaluation, measurement of KPIs, reward engineering are important

▶ MCUNet v1[9] & v2[10] are top class works - based on one shot NAS[11], Search Space reduction highest FLOPs within SRAM constraint - closely tied with TinyEngine IR for HW params

▶ μ-NAS[12] based on AE & BO, use working sets to approximate PMU & #MAC for latency

**[8] SwiftNet: Using Graph Propagation as Meta-knowledge to Search Highly Representative Neural Architectures,** Cheng *et al.*

**[9] MCUNet: Tiny Deep Learning on IoT Devices**, Ji Lin *et al.*, https://arxiv.org/pdf/2007.10319

**[10] MCUNetV2: Memory-Efficient Patch-based Inference for Tiny Deep Learning**, Ji Lin *et al.*, https://arxiv.org/abs/2110.15352

**[11] Understanding and Simplifying One-Shot Architecture Search**, Bender *et al.*, https://proceedings.mlr.press/v80/bender 18a

**[12] μNAS:** Constrained Neural Architecture Search for Microcontrollers, Edgar Liberis *et al.*, https://arxiv.org/pdf/2010.14246

tcs Research

Inventing for impact

# NAS Research & TinyML: TinyML Goals

▶ 1000-layer DNN for MCUs ? - 20-30 layer DNN fits ; Global search space good enough; complex, unexpected interconnect[8] OK ; Training smaller models marginalizes many problems

▶ For TinyML model size, peak RAM usage, #MAC are highly important

▶ Evaluation, measurement of KPIs, reward engineering are important

▶ MCUNet v1[9] & v2[10] are top class works - based on one shot NAS[11], Search Space reduction highest FLOPs within SRAM constraint - closely tied with TinyEngine IR for HW params

▶ μ-NAS[12] based on AE & BO, use working sets to approximate PMU & #MAC for latency

▶ MicroNets[13] most promising: DARTS + latency, memory, energy - all approximated by #FLOPS, MO optimization DARTS Loss + 'some regularizer' - not clear: https://github.com/liyunsheng13/micronet

▶ Moreover, this CVPR '20 paper showed that #MAC not good proxy for latency

[10] **MCUNetV2: Memory-Efficient Patch-based Inference for Tiny Deep Learning**, Ji Lin *et al.*, https://arxiv.org/abs/2110.15352

[11] **Understanding and Simplifying One-Shot Architecture Search**, Bender *et al.*, https://proceedings.mlr.press/v80/bender 18a

[12] **μNAS:** Constrained Neural Architecture Search for Microcontrollers, Edgar Liberis *et al.*, https://arxiv.org/pdf/2010.14246

[13] **MicroNets: Neural Network Architectures for Deploying TinyML Applications on Commodity Microcontrollers**, Colby Banbury *et al.*, https://arxiv.org/abs/2010.11267

[14] **Latency-Aware Differentiable Neural Architecture Search,** Yuhui Xu *et al.*, https://arxiv.org/abs/2001.06392

tcs Research

Inventing for impact

# NAS Research & TinyML: Reward Engineering

- Generally, for NAS reward is ACC - proven on benchmark datasets
- TinyML objective: Pareto Optimal - one objective can't be improved without making another worse
- argmin α ∈ A { 1.0 − ACC(α), SIZE(α), PMU(α), LAT(α) }
- Easiest method - treat all other objectives as constraints ( comes from platform)

tcs Research

# NAS Research & TinyML: Reward Engineering

▶ Generally, for NAS reward is ACC - proven on benchmark datasets

▶ TinyML objective: Pareto Optimal - one objective can't be improved without making another worse

▶ argmin α ∈ A { 1.0 − ACC(α), SIZE(α), PMU(α), LAT(α) }

▶ Easiest method - treat all other objectives as constraints ( comes from platform)

▶ DPP-Net[15] -  ACC, #params, inference time, memory usage on actual hardware

▶ MONAS[16] - scalarization: $R = \alpha * ACC − (1−\alpha) * ENERGY$

▶ µ-NAS[12] - scalarization: $L^t(\alpha) = \max(\lambda^t_1 (1.0 − ACC(\alpha)), \lambda^t_2 PMU(\alpha), \lambda^t_3 SIZE(\alpha), \lambda^t_4 MACS(\alpha))$

[15] **DPP-Net: Device-aware Progressive Search for Pareto-optimal Neural Architectures,** Jin-Dong Dong *et al.*,
https://arxiv.org/abs/1806.08198

[16] **MONAS: Multi-Objective Neural Architecture Search using Reinforcement Learning**, Chi-Hung Hsu *et al.*,
https://arxiv.org/abs/1806.10332

tcs Research

Inventing for impact

# NAS Research & TinyML: Reward Engineering

▶ Generally, for NAS reward is ACC - proven on benchmark datasets

▶ TinyML objective: Pareto Optimal - one objective can't be improved without making another worse

▶ argmin $\alpha \in A$ { $1.0 - ACC(\alpha)$, $SIZE(\alpha)$, $PMU(\alpha)$, $LAT(\alpha)$ }

▶ Easiest method - treat all other objectives as constraints ( comes from platform)

▶ DPP-Net[15] - ACC, #params, inference time, memory usage on actual hardware

▶ MONAS[16] - scalarization: $R = \alpha * ACC - (1-\alpha) * ENERGY$

▶ μ-NAS[12] - scalarization: $L^t(\alpha) = \max(\lambda^t_1 (1.0 - ACC(\alpha)), \lambda^t_2 PMU(\alpha), \lambda^t_3 SIZE(\alpha), \lambda^t_4 MACS(\alpha))$

▶ MicroNets[13] - using regularizer with DARTS: $\Sigma^K_{k=1} z_k |\theta_k|$ - how ?

**[15] DPP-Net: Device-aware Progressive Search for Pareto-optimal Neural Architectures,** Jin-Dong Dong *et al.*,
https://arxiv.org/abs/1806.08198

**[16] MONAS: Multi-Objective Neural Architecture Search using Reinforcement Learning**, Chi-Hung Hsu *et al.*,
https://arxiv.org/abs/1806.10332

tcs Research

Inventing for impact

# The Naive Neural Architecture Search: Model

- input images x $\in \mathbb{R}^{h \times w \times c}$ from an input space $X$, and a set of labels x $\in \mathbb{R}$ from an output space $Y$
- Mapping η: $X, Y, P \rightarrow \mathcal{N}(\theta)$
- $\mathcal{N}^m(\theta_\square) \circ \ldots \mathcal{N}^2(\theta_2) \circ \mathcal{N}^1(\theta_1)$, where each layer implementation $\mathcal{N}^i$ is parametrized by $\theta_i$, and implements the functionality $f_i$



- $y = f_\square( f_{\square-1}(\ldots(f_1(x))\ldots))$, where each $f_i$ can be an affine transform, non-linear transform, etc.

```
model = nn.Sequential(nn.Conv2d(1, 6, kernel_size=5, stride=1, padding=0),
                      nn.ReLU(),
                      nn.MaxPool2d(kernel_size = 2, stride = 2)),
                      nn.Linear(400, 10),
                      torch.nn.Softmax(dim=10))
```

Inventing for impact

# Naive NAS: Representing $f_\square( f_{\square-1}(\ldots(f_1(x))\ldots))$

1. Assign identifiers for each component of the architecture
2. Let the encoding of a single layer be: [ **<layer type><input><output><kernel size><stride><padding>... ]**
3. Let the architecture be a concatenation of all such layers in sequence:
   **[ [<layer 1>**<param 1><param 2>**...]**
   **[<layer 2>**<param 1><param 2>**...]**
   **[<layer n>**<param 1><param 2>**...] ]**

| Operation | Params | Values |
|-----------|--------|--------|
| Conv2D | Kernel size | {1,3,...} |
| | Out Channel | {16,32,...} |
| | In Channel | {16,32...} |
| | Stride | {1,2} |
| Linear | In | {10,64,..} |
| | Out | {10,64,..} |
| Softmax | | |

| Operation | Encoding |
|-----------|----------|
| Conv2D | 1 |
| Linear | 2 |
| MaxPool2D | 3 |
| BN | 4 |
| Softmax | 5 |

tcs Research

Inventing for impact

# Algorithm 1: Naive NAS - Random Architecture Search

**Initialize**: Set max_layers, max_params, max_score, Layers, Params, Values lists & empty architecture $\pi$

while $\rho$ < max_score do

  for each *l* in **Layers** do

    for each *h* in **Params** [ *l* ] do

      2. *l* [ *h* ] $\leftarrow$ random_choice ( **Values** [ *h* ] )

    3. Append *l* to architecture $\pi$

  4. Convert string architecture $\pi$ to $\mathcal{M}(\theta)$

  5. Train and evaluate $\mathcal{M}(\theta)$ to get *ACC, SIZ, RAM, MAC*

  6. Convert *ACC, SIZ, RAM, MAC* into a weighted score $\rho$

for each *l* in **Layers** do

  for each *h* in **Params**

    random_apply($\pi$[*l* [*h* ] ], **Values**(*h*))

tcs Research

Inventing for impact

# Enhancing Naive NAS

👎 Problems of Naive NAS
- ✖ Learns nothing, better architecture are NOT generated progressively
- ✖ We may get a good architecture, but no guarantee when

tcs Research

# Enhancing Naive NAS

👎 Problems of Naive NAS

    ✖ Learns nothing, better architecture are NOT generated progressively

    ✖ We may get a good architecture, but no guarantee when

❓ How to make the search intelligent?

    🔧 Consider and incorporate the feedback generated from the evaluation of the network

    🔧 Should strive for generating progressively better architectures

tcs Research

# Enhancing Naive NAS

👎 Problems of Naive NAS
  ✖ Learns nothing, better architecture are NOT generated progressively
  ✖ We may get a good architecture, but no guarantee when

❓ How to make the search intelligent?
  🔧 Consider and incorporate the feedback generated from the evaluation of the network
  🔧 Should strive for generating progressively better architectures
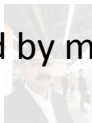
💡 Sometimes however, it works surprisingly well*

[ * ] **Random Search and Reproducibility for Neural Architecture Search,** Liam Li, Ameet Talwalkar, *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, PMLR, 2020

tcs Research

Inventing for impact

# Enhancing Naive NAS

👎 Problems of Naive NAS
  ✖ Learns nothing, better architecture are NOT generated progressively
  ✖ We may get a good architecture, but no guarantee when

❓ How to make the search intelligent?
  🔧 Consider and incorporate the feedback generated from the evaluation of the network
  🔧 Should strive for generating progressively better architectures

💡 Sometimes however, it works surprisingly well*
💡 Making the search intelligent is just putting it into one of the existing formulae - not very difficult if we know the formula & understand our problem

[ * ] **Random Search and Reproducibility for Neural Architecture Search,** Liam Li, Ameet Talwalkar,
    *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, PMLR, 2020

Inventing for impact

# Reinforcement Learning Based NAS: Model

- We model the search (NAS **agent**) as Markov Decision Process (MDP)

- Agent starts with an initial architecture $n$ - its initial **state**

- **States** are associated with $n$, identified by a state vector

- From one **state** other **states** can be generated by modifying $n$

- Modifying $n$ is taking some **action**

- There can be several choices of action from a current state

- Why are we doing all these? We want the agent to choose best **actions**

tcs Research

Inventing for impact

# RL NAS: How to Identify the Best Action?

- Agent aims to take actions, such that a *reward* value is maximized

- **Reward** represents the feedback that agent receives from *environment*

- **Reward** can be the *validation accuracy*

tcs Research

Inventing for impact

# RL NAS: How to Identify the Best Action?

✦ Agent aims to take actions, such that a *reward* value is maximized

✦ **Reward** represents the feedback that agent receives from *environment*

✦ **Reward** can be the *validation accuracy*

✦ **Reward** can be ***immediate*** or in distant ***future*** for a state trajectory

$$R_t = \rho_{t_1} + \gamma \rho_{t_2} + \gamma^2 \rho_{t_3} + \ldots$$

tcs Research

Inventing for impact

# RL NAS: Value of a State Transition

- Need to associate a *value* for transition from one state to other on an action

- Cumulative *value* of a state transition from **C□₁** to **C□₂,** on an *action* α

$$V(C_{t_1}, \alpha) = \rho^{\alpha}_{C_{t_1} C_{t_2}} + \gamma \mathbb{V}_{C_{t_2}}$$

tcs Research

# RL NAS: Value of a State Transition

➥ Need to associate a **value** for transition from one state to other on an action

➥ Cumulative **value** of a state transition from **C☐₁** to **C☐₂,** on an **action** α

$$V(C_{t_1}, \alpha) = \rho^{\alpha}_{C_{t_1} C_{t_2}} + \gamma \mathbb{V}_{C_{t_2}}$$

➥ A greedy strategy considers the maximum expected reward starting from **C☐₂**

$$V(C_{t_1}, \alpha_{t_1}) = \rho^{\alpha_{t_1}}_{C_{t_1} C_{t_2}} + \gamma \max_{\alpha_{t_2}} V(C_{t_2}, \alpha_{t_2})$$

➥ **V(∘) → action-value function** and the individual **V(C☐, α) → Q-values**

tcs Research

Inventing for impact

# Algorithm 2: RL NAS

**Initialize**: Set max_score, Layers, Params, Values lists, empty architecture $n$, action set $\mathcal{A}$ & state set $\mathcal{C}$

while $\rho$ < max_score do

  for each $l$ in **Layers** do

     for each $h$ in **Params** [ $l$ ] do

       2. $l$ [ $h$ ] $\leftarrow$ random_choice ( **Values** [ $h$ ] )

     3. Append $l$ to architecture $n$

     4. Update $\mathcal{A}$ with $\alpha$, $\mathcal{C}$ with $C_n$

               **EXPLORE**

4. Convert string architecture $n$ to $\mathcal{M}(\theta)$

5. Train and evaluate $\mathcal{M}(\theta)$ to get *ACC, SIZ, RAM, MAC*

6. Convert *ACC, SIZ, RAM, MAC* into a weighted score $\rho$

               **EVALUATE**

  for each $C_n$ in $\mathcal{C}$ do

    for each $\alpha$ in $\mathcal{A}$ do

    7. Use $C_n$, $\alpha$, $\rho$, $V(C_n, \alpha)$ to train $\mathcal{M}(\theta)$

# Algorithm 2: RL NAS

**Initialize**: Set max_layers, max_params, max_score, Layers, Params, Values lists, empty architecture $\eta$, & action set $\mathcal{A}$

while $\rho$ < max_score do

  If *random_value* < $\varepsilon$ do

    1. **EXPLORE**

  else do

       for each *l* in **Layers** do

        2. $C_{\eta+1} \leftarrow \mathcal{M}(\theta)\,(C_\eta)$

        3. $l \leftarrow C_{\eta+1}$

        4. Append *l* to architecture $\eta$

**EXPLOIT**

  7. **EVALUATE**

# Case Study: RL-NAS

**Generating Tiny Deep Neural Networks for ECG Classification on MCU**, Shalini Mukhopadhyay, Swarnava Dey, Avik Ghose (TCS Research), Pragya Singh (IIIT-D), Pallab Dasgupta (IIT KGP) IEEE Percom 2023

Inventing for impact

# Case Study: RL-NAS

- RL based NAS (MetaQNN[2]) - ECG Atrial Fibrillation, Smoking episode, & Hand Gesture
- 3000 iterations, 90 hours (A100), 20-layer, 225 KB model with SOTA F1 score for ECG-AF
- MO - ACC, SIZE, SRAM, #MAC: $R = (W_a ACC + \sum_i W \exp(w_i P_i)) / \sum |W|$



**Multi-objective Convergence**



**Clear Front of good models**



Adafruit Feather
NRF52840
RAM Size : 256 kB
Flash Size : 1 MB
Target MACs : 16 M

**Search Parameters**

**Generating Tiny Deep Neural Networks for ECG Classification on Micro-controllers**, Mukhopadhyay et al., IEEE Percom 2023

tcs Research

Inventing for impact

# Algorithm 3: EVO NAS

**Initialize**: Set max_score, max_pop, Layers, Params, Values lists, empty architecture $n$, action set $\mathcal{A}$ & state set $\mathcal{C}$

1. $POPULATION \leftarrow \{\phi\}$

for 1… max_score do

  2. $n, \mathcal{A} \leftarrow$ **EXPLORE**

  3. $\rho \leftarrow$ **EVALUATE**($n$)

  4. $POPULATION$.append($n, \mathcal{A}, \rho$)

$mut\_rate \leftarrow 1.0$

for 1… $max\_score$ do

  5. $parent \leftarrow POPULATION.fittest()$

  6. $child \leftarrow$ **mutate**($parent, mut\_rate$)

(contd.)

for 1… $max\_score$ do

  7. $parent \leftarrow POPULATION.fittest()$

  8. $Child, \mathcal{A} \leftarrow$ **mutate**($parent, mut\_rate$)

  9. $\rho \leftarrow$ **EVALUATE**($child$)

  10. $POPULATION$.append($child, \mathcal{A}, \rho$)

  11. $POPULATION.remove\_least\_fits()$

  12. Update $mut\_rate$

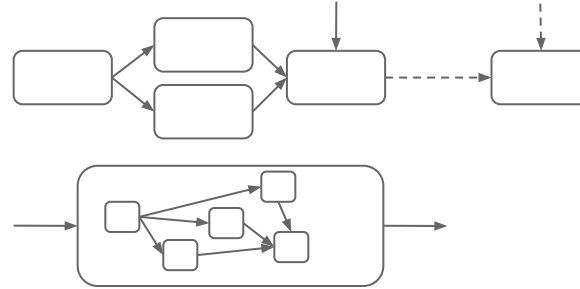13. $final\_architectures \leftarrow POPULATION.fittest()$

# Points to Note

✋ The main approach is sampling architectures randomly, evaluating, learning to sample better architectures in an informed manner

✋ Sequentially building network

✋ Reward is completely decoupled from search

✋ Discrete choices for layers and params

tcs Research

Inventing for impact

# DARTS: Differentiable NAS

Considers a fixed network a number of *cells*

*Cells* have nodes, connectivity (ops) is searched

Research

# DARTS: Differentiable NAS

Considers a fixed network a number of *cells*

*Cells* have nodes, connectivity (ops) is searched



training →

$x^{(j)} = \sum_{i<j} o^{(i,j)}(x^{(i)})$

**(a)**

training →

$x^{(j)} = \sum_{i<j} o^{(i,j)}(x^{(i)})$

**(b)**

training →

$x^{(j)} = \sum_{i<j} o^{(i,j)}(x^{(i)})$

**(c)**

$$\bar{o}^{(i,j)}(x) = \sum_{o \in \mathcal{O}} \left( \left( exp(\alpha^{(i,j)}_o) \right) / \left( \sum_{o' \in \mathcal{O}} exp(\alpha^{(i,j)}_{o'}) \right) \right) * o(x)$$

# Algorithm 4: DARTS – Differentiable Architecture Search

Set max_layers, max_params, Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge **(i, j)**

while not converged do

      1. Update architecture **α** by descending $\nabla\alpha L_{val}(w - \xi \nabla w \, L_{train}(w, \alpha), \alpha)$

      ($\xi = 0$ if using first-order approximation)

      2. Update weights **w** by descending $\nabla w L_{train}(w, \alpha)$

3. Discretize for the final architecture based on the learned **α**.

4. Re-train

$$min_{\alpha}L_{val}(\alpha, argmin_{w}L_{train}(\alpha, w)) \approx \nabla\alpha L_{val}(w - \xi \nabla w \, L_{train}(\alpha, w)$$

# Take Home Points

✋ Gordon Moore is no more. Moore's law is supposed to reach saturation levels.
   Still, it seems that the capacity of tiny devices will keep increasing in this decade

✋ Consequently, we have to search larger & complex networks for *Edge* devices

✋ Thus multi-objective, one-shot differential NAS is *the* research direction for TinyML

✋ Meanwhile, this tutorial showed that the existing multi-objective, sample-based NAS methods are <u>quite usable</u> for the next few years

tcs Research

Inventing for impact

*Thank You for bearing with me!* 😊

*Questions* 👂

Research

Inventing for impact

# Copyright Notice

www.tinyml.org

# Copyright Notice

## www.tinyml.org