

This work has received funding from the European Union's Horizon research and innovation programme under grant agreement No 101070374.



Towards a tailored mixed-precision sub-8-bit quantization scheme for Gated Recurrent Units using Genetic Algorithms

Riccardo Miccini^{1,2}, Alessandro Cerioli^{1,2}, Clément Laroche¹, Tobias Piechowiak¹, Jens Sparsø², and Luca Pezzarossa²

¹GN Store Nord / ²Technical University of Denmark

TinyML Research Symposium

April 22, 2024

Motivations and Contributions

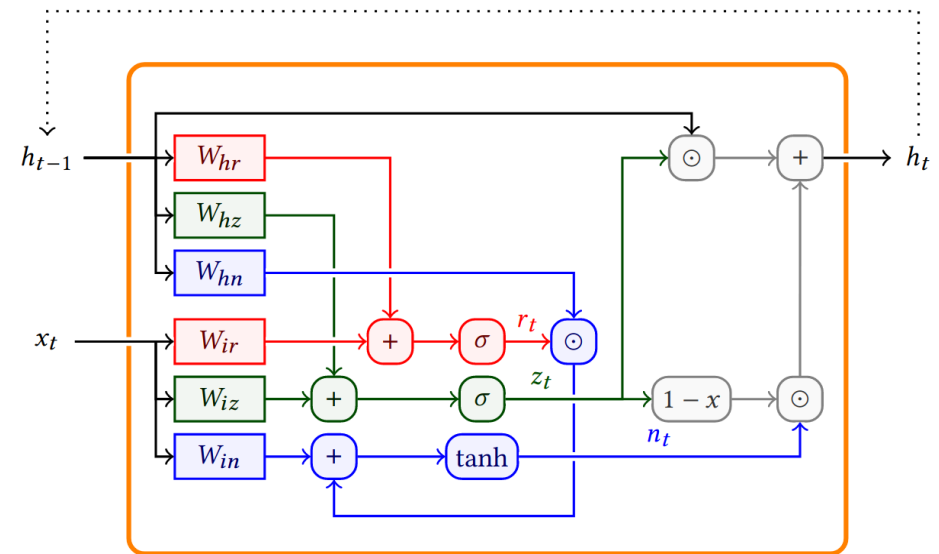
- AI for audio processing on embedded devices
 - Sub-8-bit arithmetic can be quite efficient
 - Existing quantization tools have **limited support for sub-8-bit or GRUs**
 - Quantization parameters can be tedious to tune
- We propose:
 - Integer modular quantization scheme for GRUs
 - Genetic search for Pareto-optimal **any-bit quantization scheme**
 - Evaluation on sequence classification tasks of varying complexity

Outline

- Background
- Methods
- Results
- Conclusion

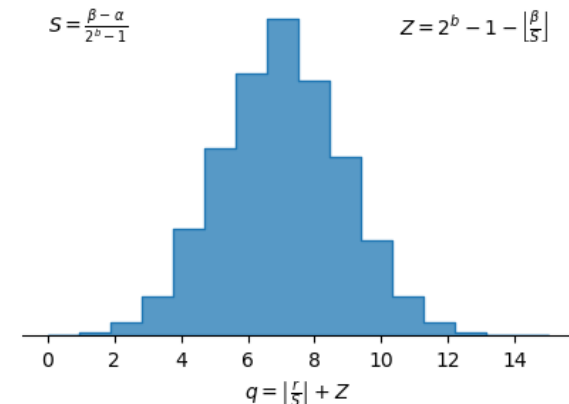
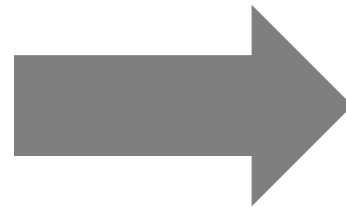
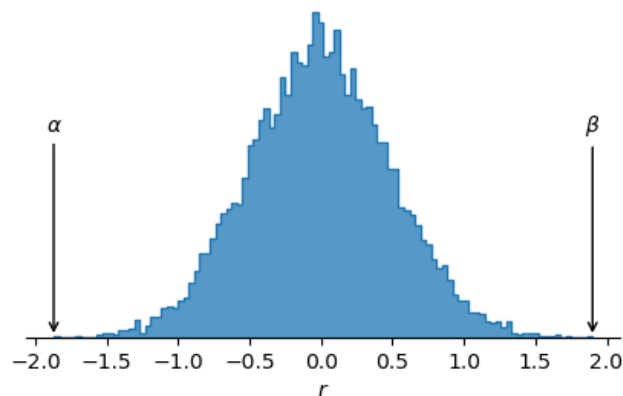
Background | Gated Recurrent Unit (GRU)

- GRUs are a class of Recurrent Neural Networks
 - Output depends on current input and previous output
 - Can **process sequential data** of arbitrary length
 - Often used in audio applications
- Like LSTMs, but more compact
 - Address vanishing/exploding gradient
 - 3 gates: **reset**, **update**, **new state**
 - Internal state is the output



Background | Quantization

- Map **continuous** data (high bit-width) **into** a **discrete** set (low bit-width)
 - Decrease storage, memory, and compute requirements
 - Lower precision might affect performances
- Linear/uniform quantization
 - Straightforward and most popular
 - Parameterized by bit-width, scaling factor S , and zero-point Z



Background | Genetic Algorithms (GA)

- Class of **optimization** algorithms
 - Inspired by biological principles (natural selection, evolution)
 - Suitable for non-differentiable multidimensional problems
- Start with a random population
 - Each specimen is a potential solution characterized by its **genome**
 - Solutions are **evaluated, ranked, and selected** for mating
 - New specimens are generated through **crossover** and **mutation**
 - Repeat with new specimens

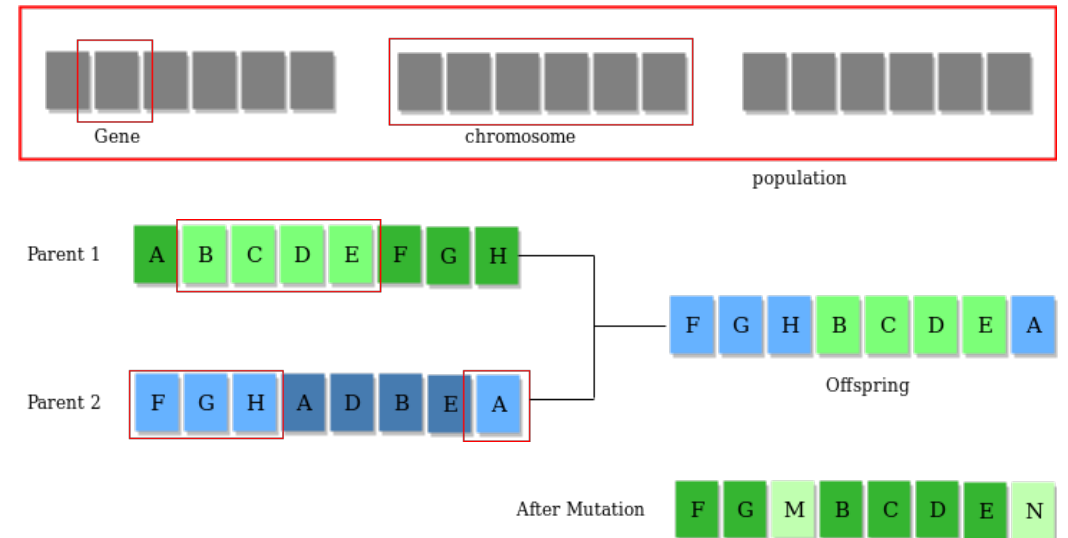
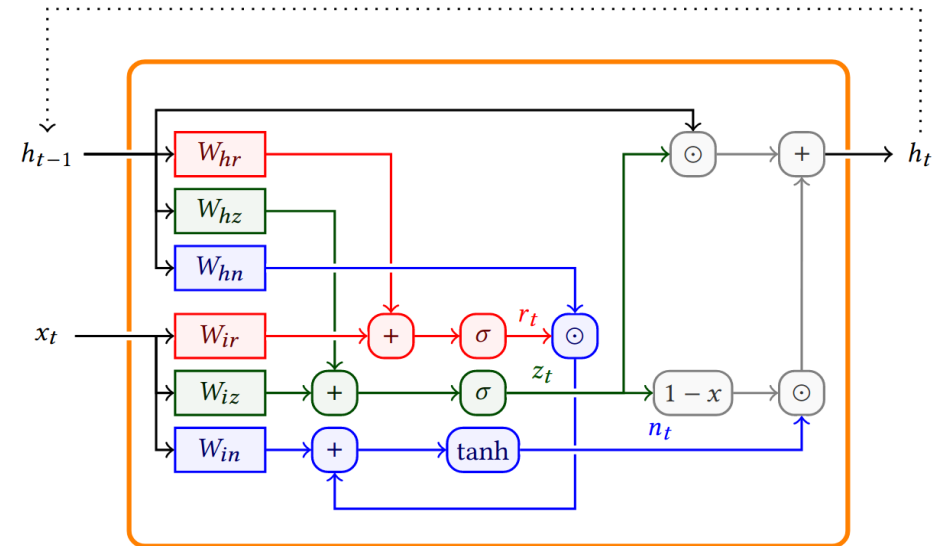


Image source: <https://www.geeksforgeeks.org/genetic-algorithms>

Methods | Modular quantization for GRUs

- Operators inside a GRU Cell
 - Linear/dense (matmul)
 - Hadamard (element-wise) product
 - Element-wise sum
 - Non-linearities: sigmoid, tanh
- Derive quantized versions
 - Substitute dequantization equation to the operator definition
 - Solve for quantized output
 - Combine scaling factors and convert to fixed-point
- For non-linearities, use look-up tables

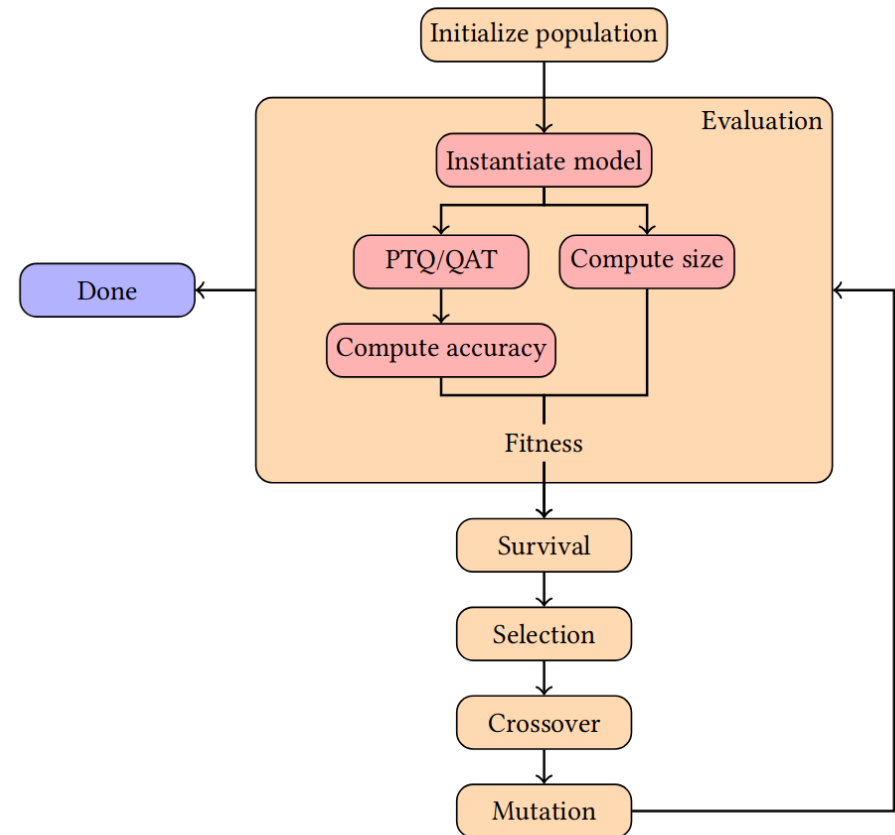


$$\begin{aligned} \mathbf{y} &= \mathbf{x}_1 + \mathbf{x}_2 \\ S_y(q_y - Z_y) &= S_1(q_1 - Z_1) + S_2(q_2 - Z_2) \\ q_y &= \frac{S_1}{S_y}(q_1 - Z_1) + \frac{S_2}{S_1}(q_2 - Z_2) + Z_y \end{aligned}$$

Methods | Quantization scheme search

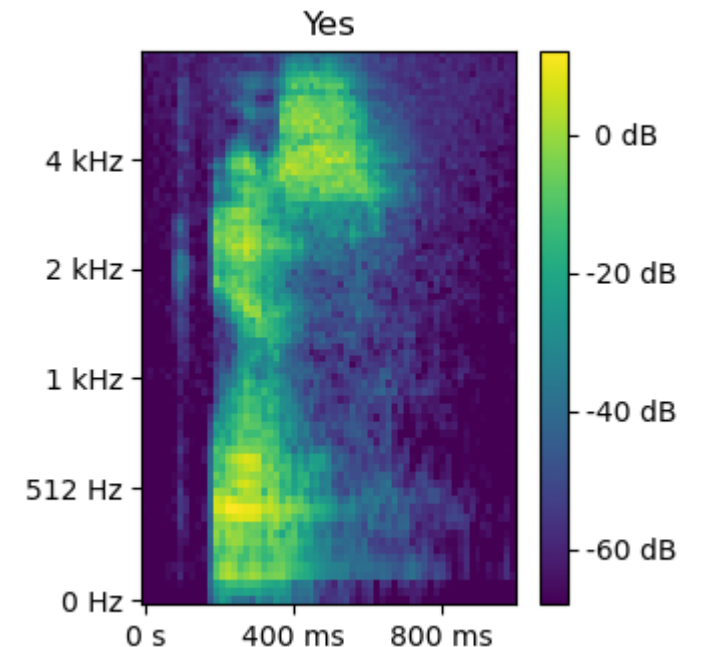
- Genetic algorithm: NSGA-II
 - Selection step based on **Pareto efficiency**
 - Genome: bit-width (from 2 to 8 incl.) of each quantized operator, 17 in total
 - 40 initial specimens, 20 generations
- Fitness metrics
 - Model **accuracy**, computed on a validation set after quantization
 - Normalized complement of **model size**, based on number of trainable parameters

$$\hat{M}^c = 1 - \frac{M_Q}{M_{FP16}} \quad \leftarrow \begin{array}{l} \text{model size with current quant. scheme} \\ \text{model size with 16-bit weights} \end{array}$$

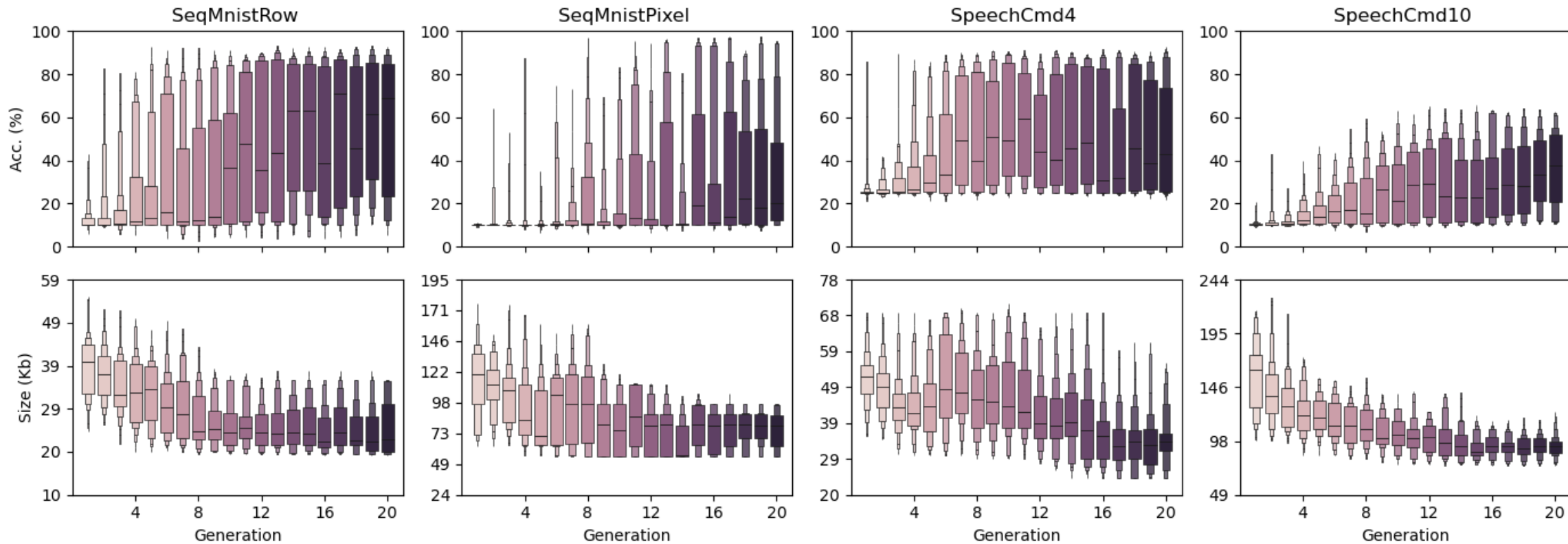


Results | Experimental setup

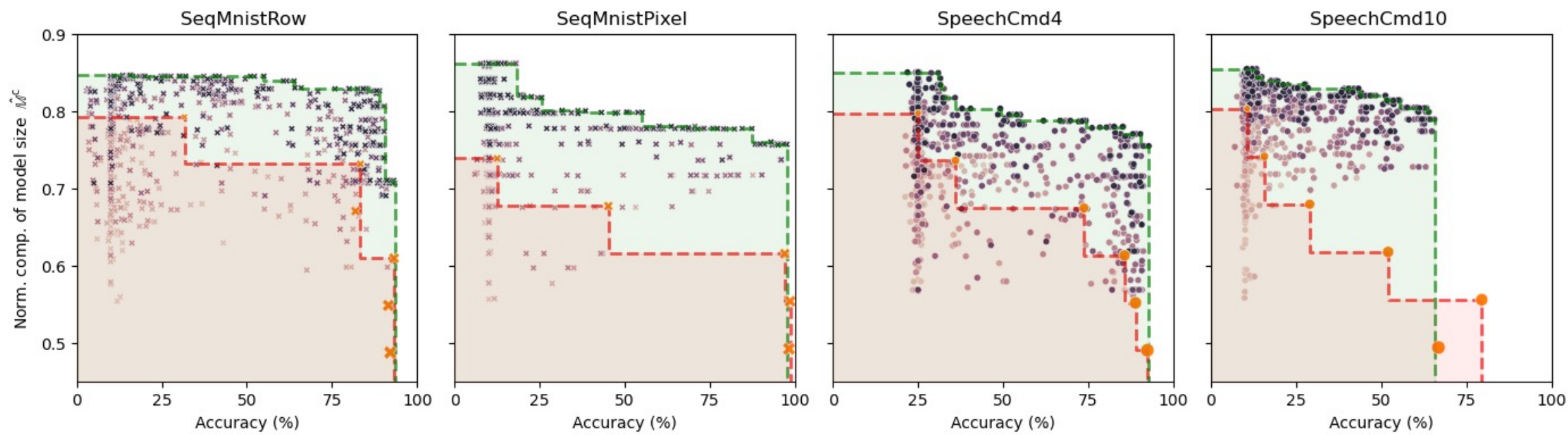
- Tasks and Datasets
 - MNIST (row-wise or pixel-wise): 28 timesteps \times 28 features or 361 scalar timesteps
 - SpeechCmd (4 and 10 keywords): 63 timesteps \times 40 features (Mel-Spectrogram)
- Model
 - GRU \rightarrow FC \rightarrow SoftMax; 256 or 128 hidden units
- Quantization
 - QAT on MNIST tasks (10% train set, learning rate = $1e-5$)
 - PTQ on SpeechCmd tasks (100% train set for calibration)
- Baseline
 - Homogeneous quantization (bit-widths from 3 to 8 incl.)
- Evaluation metrics
 - Same as fitness metrics



Results | Genetic search



Results | Pareto fronts



Generation	Homogeneous Bits	Quantization Mode
0	3	● ptq
4	4	× qat
8	5	
12	6	
16	7	
Baseline	8	
		Pareto Fronts
		--- homogeneous
		--- mixed-precision

Results | Observations

- Over the generations in the search:
 - Average accuracy increases
 - Model size decreases
- When using genetic search, we:
 - Exceed the Pareto fronts of baselines
 - Maintain 8-bit baseline performances
 - Achieve a model size reduction between 25% and 55%
- ...except for SpeechCmd10
 - We experience a 17% drop in max accuracy
 - High bit-width solutions are not explored

Conclusion | Limitations and Future work

- Most of the derived solutions are low in accuracy
 - Constrain fitness metrics during survival step
 - Split search into exploration and exploitation phases
- Can it scale to more challenging scenarios?
 - Apply to larger tasks/models (e.g. speech enhancement)
 - Extend genome (e.g. observer types and granularity)
- Need for dedicated hardware support
 - Some SoCs might support this (e.g. GAP9 vectorized 2/4-bit arithmetic)

Conclusion | Summary

- In this work, we:
 - Introduce a modular integer quantization scheme for GRUs
 - Apply a multi-objective genetic search of quantization parameters
 - Evaluate the system on sequence classification tasks
- Results:
 - Heterogeneous solutions are more Pareto-efficient
 - Substantial decrease in model size with comparable accuracy

Copyright Notice

This presentation in this publication was presented at the tinyML[®] Research Symposium 2024. The content reflects the opinion of the author(s) and their respective companies. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

www.tinyml.org