

Free Bits: Latency Optimization of Mixed-Precision Quantized Neural Networks on the Edge

Georg Rutishauser¹ Prof. Dr. Francesco Conti² Prof. Dr. Luca Benini^{1,2} georgr@iis.ee.ethz.ch f.conti@unibo.it lbenini@iis.ee.ethz.ch

@pulp_platform 🔰 pulp-platform.org 🎻

¹Integrated Systems Laboratory, ETH Zürich ²Università di Bologna

youtube.com/pulp_platform

Edge AI: Quantization for Energy Efficiency

- Insatiable hunger for data number of IoT (sensor) nodes in use exploding:
 83 Bn. devices by 2024! [1]
- On-device processing of collected data: Edge AI
- ML models tend to be large and compute-intensive
 - Particularly for difficult tasks



2

Quantization Isn't Free Lunch!



^[2] Esser et al., "Learned Step Size Quantization" y, 2023

Fine-Grained Quantization: Mixed-Precision Networks



- Each weight and activation tensor is quantized to a different precision
- Maximize improvement in target metric (latency, model size, memory footprint), minimize accuracy drop
- But: Exponential number of configurations (MNv1: 2×10^{25})

Mixed-Precision Search:

- Goal: map Pareto front of accuracy-cost trade-off
 - Real-world cost: latency, model size, memory footprint...



[1] M. van Baalen et al., Bayesian Bits: Unifying Quantization and Pruning

BOPs Don't Correspond to Inference Latency

- Let's find & deploy a mixed-precision network to mixed-precision hardware!
 - Algorithm: Bayesian Bits



On real systems, latency is not proportional to BOPs!

The BOP-Latency Gap

Profiling ALL MNv1 layers in ALL supported precisions (~1-2h):



Counterintuitive: some layers run faster in higher precisions!

Free Bits

- Increasing their precision ("Free Bits") gives us:
 - Lower Latency
 - Higher Accuracy
- Setup:
 - Assume layer-wise execution $\rightarrow L_{net} \approx \sum_l L_l$
 - Given network architecture and target platform, find latency-optimized MP configurations

FREE BITS: Increase a layer's precision if it decreases latency

Experiments

- Networks: MobileNetV1-224-0.75, MobileNetV2-224-1.0
- Task: ImageNet, 1000-class classification
- Target platforms: PULP systems with different ISA Extensions
 - Parallel Ultra Low Power: Multicore RISC-V MCUs with specialized DSP ISA extensions
 - Computation on **cluster** of 8 RISC-V cores
 - 512 KiB main L2 memory, 64 KiB cluster L1 memory
 - Networks mapped with DORY [1] and run on a cycle-accurate simulator (GVSOC [2])
 - 3 RISC-V ISA Extensions (2 shown in this presentation):
 - XPulpV2 SIMD 8-bit arithmetic, post-increment load/store, ...
 - XPulpNN XPulpV2 + SIMD sub-byte operations on equal-precision operands
 - 2 variants of Bayesian Bits
 - Vanilla as in original paper
 - Locked weight and activation precisions constrained to be equal

Results: MobileNetV1 on XPulpNN



MobileNetV1 on Different ISA Extensions

	XPulpNNV1: SW Unpacking		XPulpV2: 8-bit SIMD only		
Network	Lat. vs. 8b	Accuracy	Lat. vs. 8b	Accuracy	
Baseline: 8b	+0%	69.1%	+0%	69.1%	
Baseline: 4b					
4b + Free Bits					
<0.5 pp. acc. drop					
	Even 8-bit-only ISA profits				
	from mixed precision!				

Conclusion: Use Your Empirical Platform Knowledge!

- BOPs are not a sufficient proxy for complexity on real platforms
- Free Bits:
 - Start with DNAS-based, HW-agnostic precision search
 - Profiling-based heuristic: Increase precision where it benefits latency
 - Improves accuracy and latency of MP networks on PULP systems
 - Applicable to **any platform**
- Benchmarks: MobileNetV1 and MobileNetV2 on PULP

✓-28% latency @ FP32 accuracy for MobileNetV1 w/ XPulpNN ISA
 ✓ Latency improvement on 8-bit only ISA (XPulpV2)

✓ Homogeneously quantized networks (8b, 4b) dominated

PULP Platform Open Source Hardware, the way it should be!

Georg Rutishauser Prof. Dr. Francesco Conti Prof. Dr. Luca Benini

georgr@iis.ee.ethz.ch f.conti@unibo.it lbenini@iis.ee.ethz.ch

Institut für Integrierte Systeme – ETH Zürich

DEI – Universitá di Bologna

ETHZÜRICH DAMA MATER STUDIORUM



@pulp_platform 🔰



youtube.com/pulp_platform

Unpacking Mismatched Operands

- Packed-SIMD operations on operands of different precisions: Lowerprecision operand must be "unpacked" to to the higher precision
- Example: MAC of 8-bit and 4-bit operands



BOPs Don't Correspond to Inference Latency

- Our target platforms: PULP systems
 - Parallel Ultra Low Power multi-core, near-threshold MCUs
 - Computation on cluster: 8×32 -bit RISC-V cores with custom ISA Extensions
 - **XPulpNN:** vectorized sub-byte MAC operations $16 \times 2b/8 \times 4b/4 \times 8b$ MACs/cycle/core
 - Let's try to find a good mixed-precision MobileNetV1 with Bayesian Bits!



⁸ May, 2023

Results: MobileNetV2 on XPulpNN



 ✓ Free Bits dominates baselines
 ✓ -11% latency @ -0.3 pp. accuracy

Q: Why is the 4b/4b baseline so slow?
A: Outputs of adder nodes are always 8b
→ precision mismatch in next layer

Full Results: MNv1, MNv2/XPV2, XPNNv1, XPNNv2

Acc. Margin	ISA	MobileNetV1		MobileNetV2	
		Lat. vs. 8b	Acc.	Lat. vs. 8b	Acc.
8b Baseline	all	+0%	69.1%	+0%	71.5%
0.5 pp.	XPv2	-5.5%	69.3%	-3.4%	71.0%
	XPNNv1	-27.9%	68.6%	-10.9%	71.2%
	XPNNv2	-28.6%	68.6%	-15.3%	71.0%
1.5 pp.	XPv2	-5.5%	69.3%	-6.3%	70.7%
	XPNNv1	-34.4%	67.6%	-15.1%	70.4%
	XPNNv2	-35.1%	67.6%	-15.3%	71.0%
4b + FB	XPv2	-3.5%	67.7%	-7.7%	70.9%
	XPNNv1	-37.1%	66.3%	-12.8%	69.9%
	XPNNv2	-39.8%	66.6%	-25.7%	69.6%
4b Baseline	XPv2	+49.9%	65.6%	+37.0%	69.3%
	XPNNv1	-32.3%	65.6%	+48.9%	69.3%
	XPNNv2	-38.3%	65.6%	-23.4%	69.3%

Integer Quantization can Help...



- Represent weights & activations as low-bitwidth integers
- Run inference using integer arithmetic only
- Benefits:
 - Smaller model size ——— Less storage needed
 - Smaller activation size —— Less memory needed
 - More ops/cycle Lower inference latency





The BOP-Latency Gap

• Profiling ALL MNv1 layers in ALL supported precisions (~1-2h):



- On real systems, latency is not proportional to BOPs!
 - HW constraints: equal-precision operands required
 - Overhead to unpack lower-precision operands
 - E.g. NVidia INT4 Tensor cores, XPulpNN
 - SW non-idealities:
 - Some kernels may see limited benefit from low-precision operands
 - Performance depends on tiling strategy 8 May, 2023

Edge AI: Energy Efficiency is Everything

- Insatiable hunger for data number of IoT (sensor) nodes in use exploding: 83 Bn. devices by 2024! [1]
- On-device processing of collected data: Edge AI
- This scaling can only be sustained if nodes become even more:
 - Cheap \rightarrow MCU-class systems
 - Powerful \rightarrow Optimization for ML algorithms
 - Versatile \rightarrow battery-powered \rightarrow energy efficiency is everything!
- ML models tend to be large and compute-intensive
 - Particularly for difficult tasks

ML
opt