

tinyML[®] Research Symposium

Enabling Ultra-low Power Machine Learning at the Edge

March 27, 2023



www.tinyML.org

Fused Depthwise Tiling for Memory Optimization in TinyML Deep Neural Network Inference

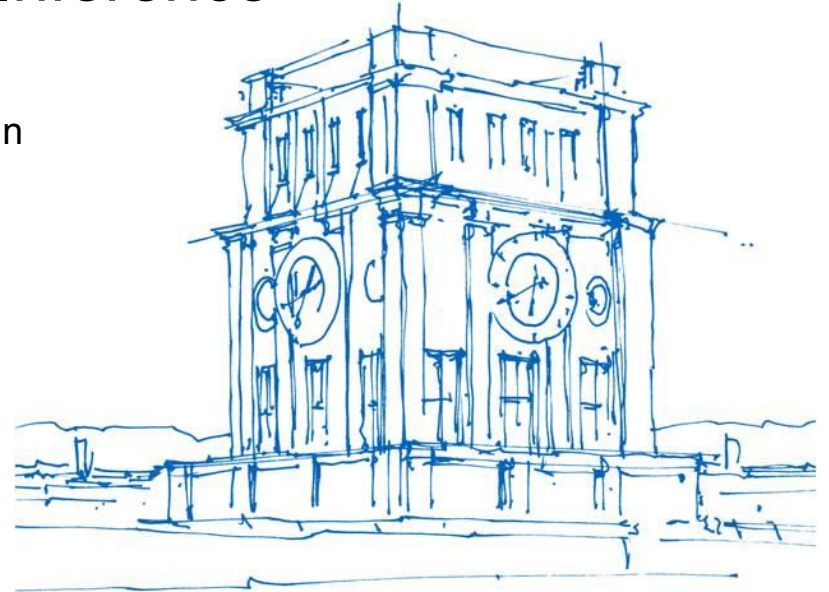
Rafael Stahl, Daniel Mueller-Gritschneider, Ulf Schlichtmann

Technical University of Munich

School for Computation, Information and Technology

Chair for Electronic Design Automation

Burlingame, 27th of March 2023




Uhrenturm der TUM

Machine Learning on Edge Devices

- Focus: **Inference**
- Improves:
 - Communication demand
 - Latency
 - Data privacy
- Many application suitable for extreme low-power: **tinyML**
 - Keyword Spotting
 - Visual Wake-up
 - Anomaly Detection
 - Radar Gesture Detection

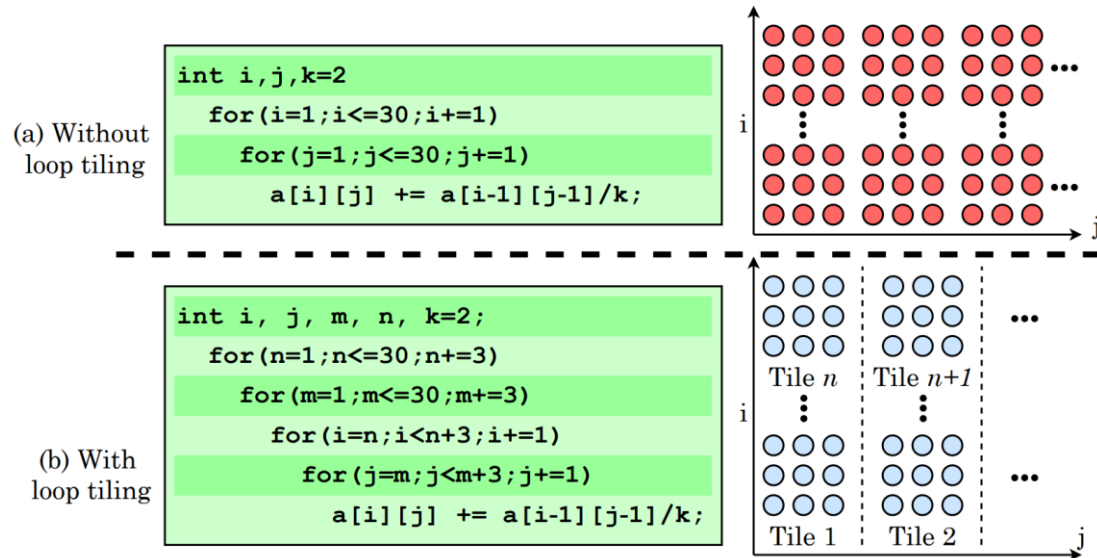


Challenge: Memory

- Power usage
 - Cost
- 
- Memory usage**
- Reducing memory with accuracy trade-off:
 - Quantization
 - Pruning
 - Network Architecture Search (NAS)

Loop Tiling

- Loop transformation to exploit spatial and temporal locality
- Typically employed for performance optimization



Fused Tiling

Tiling:

- Compute large intermediate buffer in multiple tiles

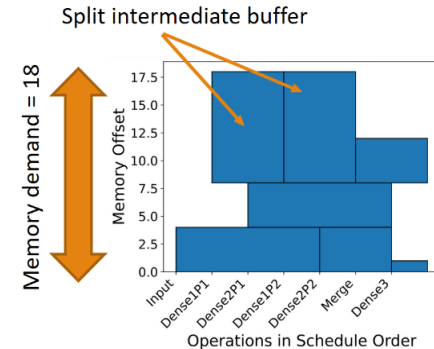
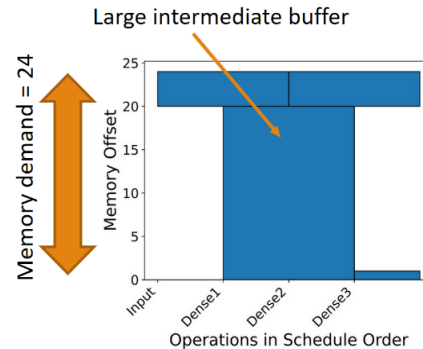
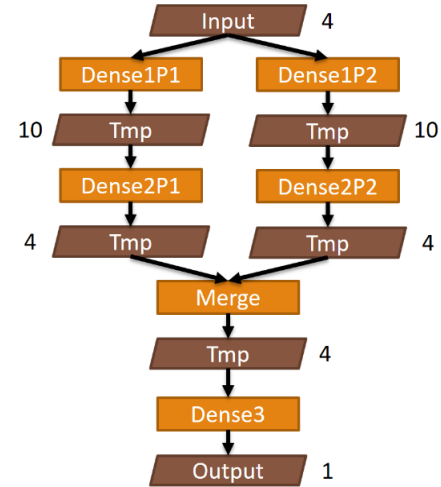
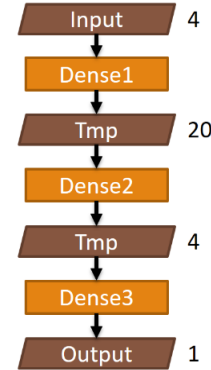
Fusion:

- Operator fusion decouples their lifetime

→ Lifetimes of split large buffers do not overlap

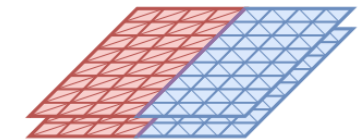
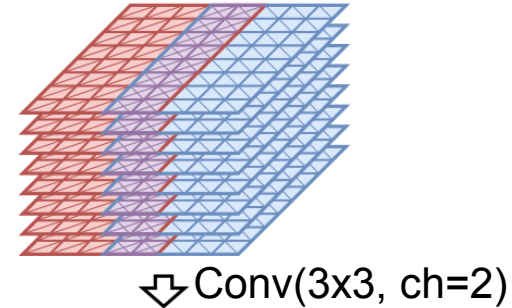
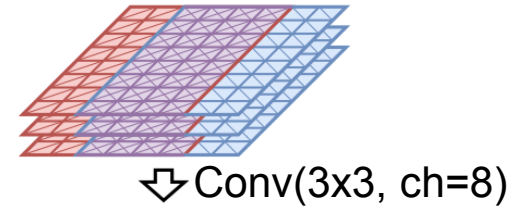
→ Their storage buffers may overlap

→ Memory reduction



Fused Feature Map Tiling (FFMT)

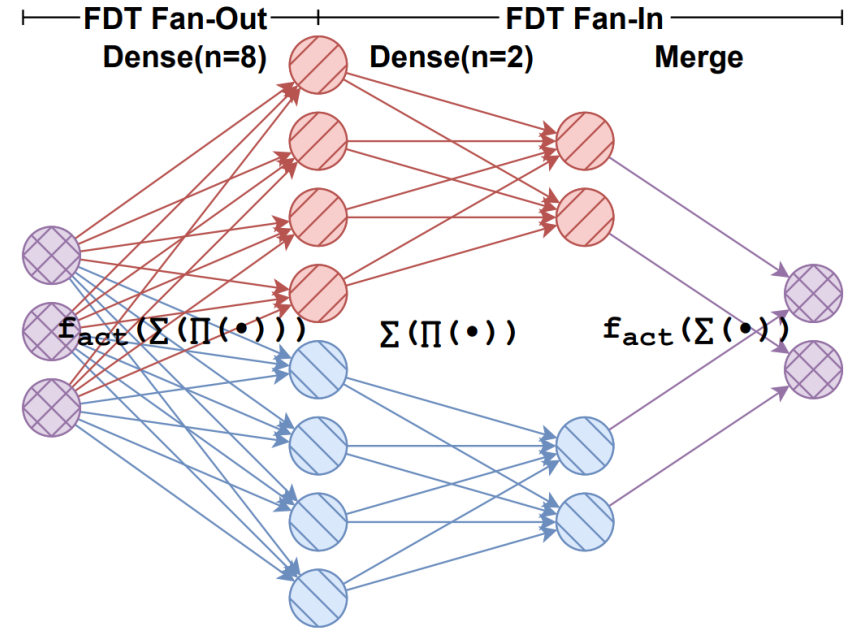
- Operator fusion through spatial locality of convolution
- Introduces overlap from kernel size
- Does not support operations with large input dependencies
 - Fully connected
 - Convolutions with very large kernel sizes



Fused **Depthwise** Tiling (FDT)

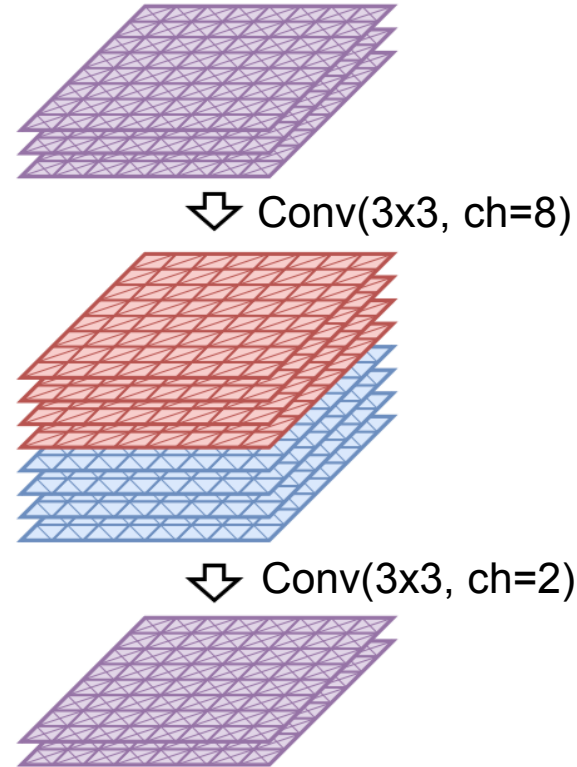
- Allows fusion of two operations with large input dependencies
- Accumulates **partial sum** in second output
- Requires **Merge** operation

- New tiling opportunities
- No significant run time overheads



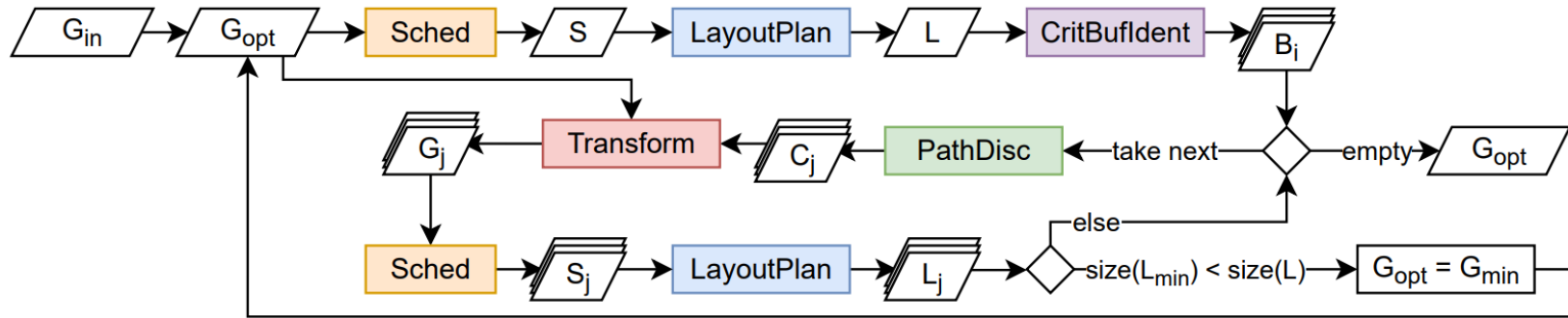
FDT for Convolutions

Split by feature maps instead of neurons

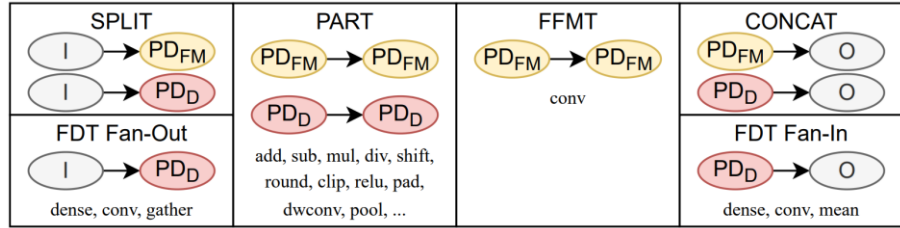


End-to-end Deployment Flow

- Determines where, and how to apply fused tiling
- Memory-aware scheduling
- Memory buffer layout planning
- Path discovery

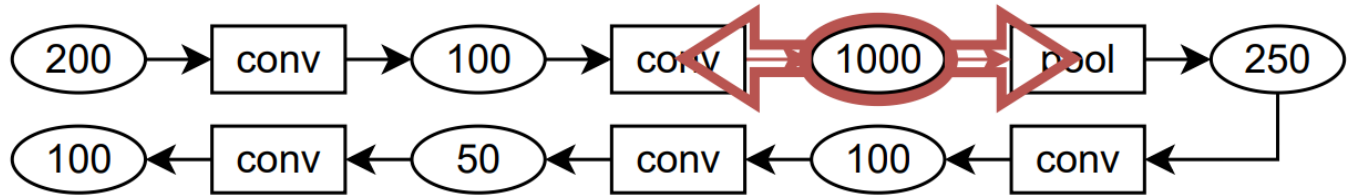


Path Discovery

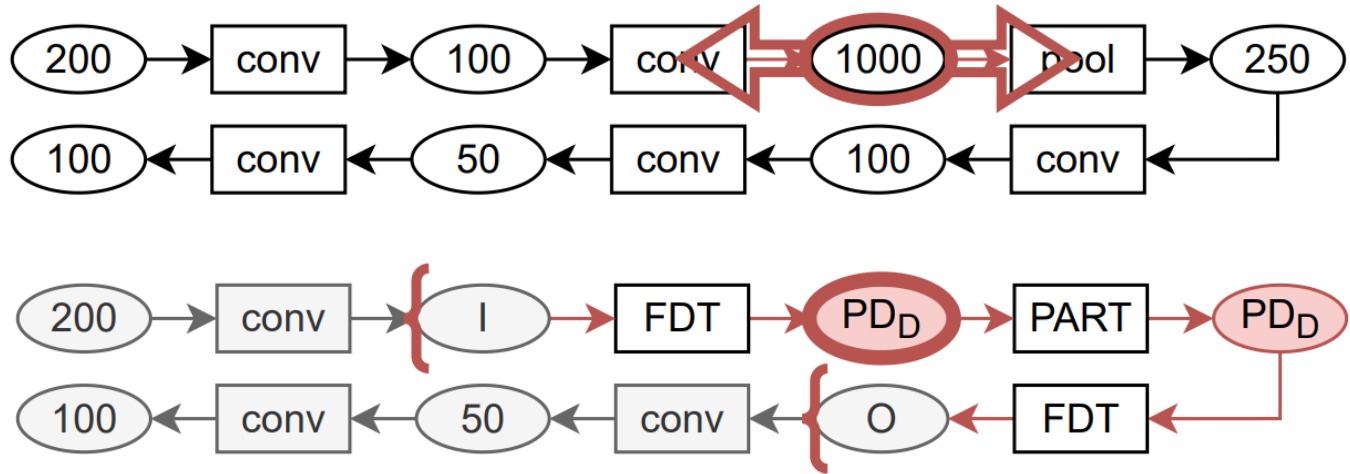
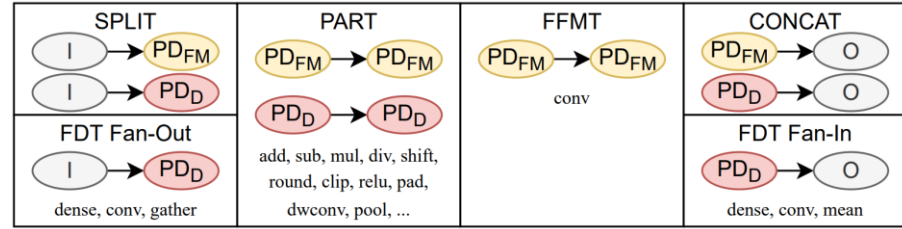


- Finds an optimized sequence of operations (**path**) for fused tiling
- Intermediate buffers and operations are matched according to **blocks**

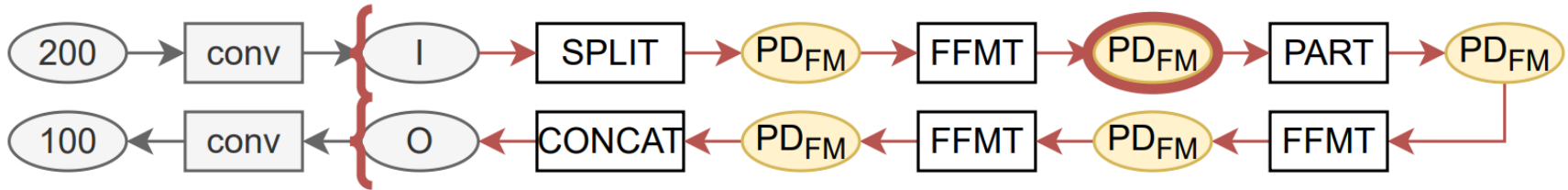
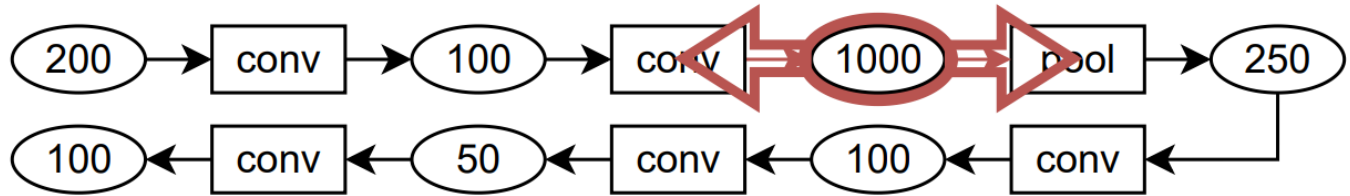
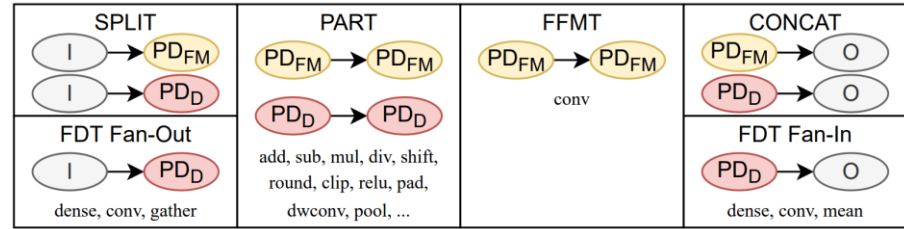
Untiled



Path Discovery – FDT



Path Discovery – FFMT



Implementation

- Implemented in Apache TVM
 - Suitable for complex transformation passes
- Evaluated seven quantized models in RISC-V RV32GC
 - Memory usage from sections of compiled binary
 - Performance estimation from multiply ops of optimized graph
- To appear open source at: <https://github.com/tum-ei-eda/moiopt>



Results

Model	Mem [kB]		[%]		MACs [1 million]			[%]		
	Untiled	FFMT	FDT	FFMT Savings	FDT Savings	Untiled	FFMT	FDT	FFMT Overhead	FDT Overhead
KWS	65.6	65.6	53.7	0.0	18.1	2.66	2.66	2.66	0.0	0.0
TXT	18.6	18.6	4.43	0.0	76.2	0.00	0.00	0.00	0.0	0.0
MW	17.6	7.04	11.3	60.9	35.5	0.06	0.06	0.06	0.0	0.0
POS	9.35k	5.11k	8.94k	45.3	4.4	837	1215	837	45.1	0.0
SSD	14.3k	8.66k	12.2k	39.4	14.6	313	314	313	0.2	0.0
CIF	179	76.7	170	57.1	5.0	5.52	6.02	5.52	9.0	0.0
RAD	36.2	26.7	29.4	26.3	18.8	0.09	0.09	0.09	0.0	0.0
Avg.				32.7	24.7				7.8	0.0

Summary

- Applied **Fused Depthwise Tiling** to DNN graphs for memory optimization
 - Built end-to-end deployment flow for evaluation
- Reduces memory usage where previously not possible
- Adds alternative solution where existing tiling causes too much performance overhead



Copyright Notice

This presentation in this publication was presented at the tinyML[®] Research Symposium (March 27,2023). 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