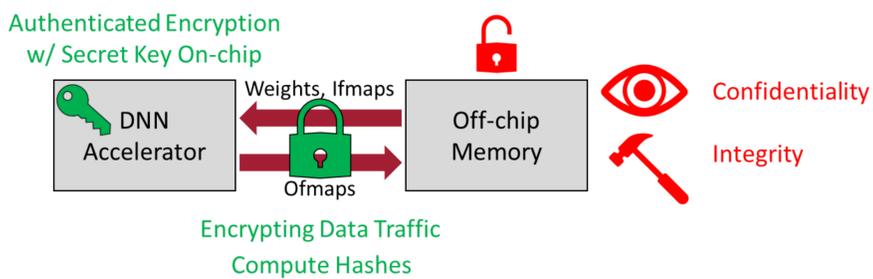
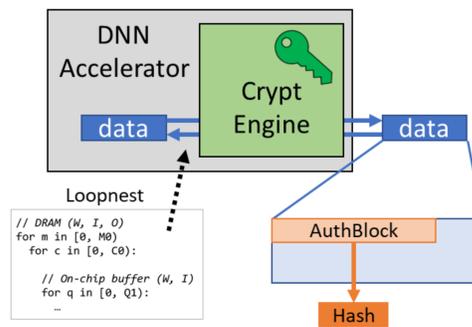


Off-chip Memory Security

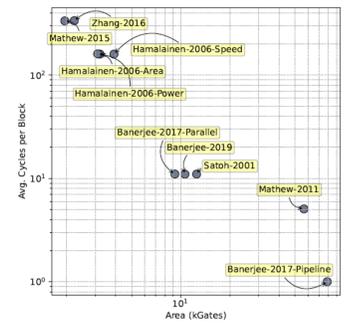


- Privacy and integrity of the untrusted off-chip memory
- Trusted execution environment (TEE): encrypt and authenticate every off-chip data traffic

Design Space Exploration Issues



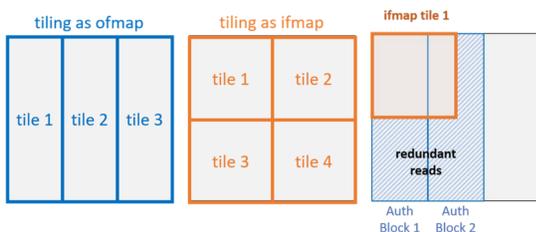
Mapping Problem
Extend to consider cryptographic operations



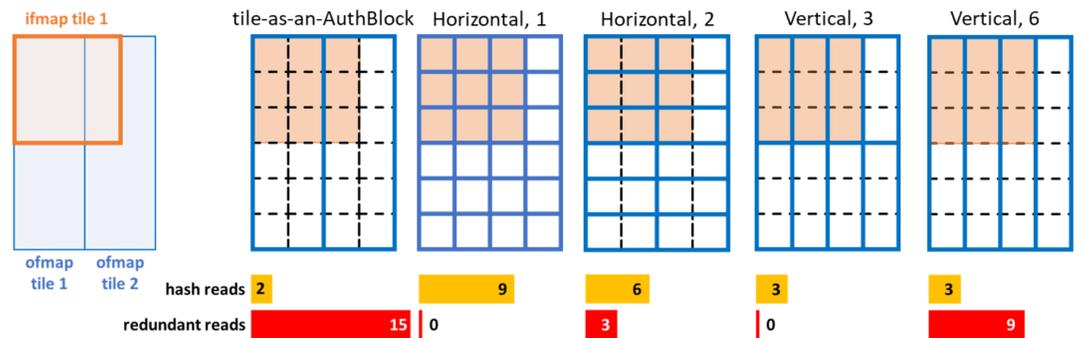
CryptEngine Overhead
Trade-off space for area, performance, energy

Scheduling Search Engine

Authentication Block Assignment



Cross-layer dependency can cause significant additional off-chip traffic when AuthBlock assignment is simply done according to the tiles



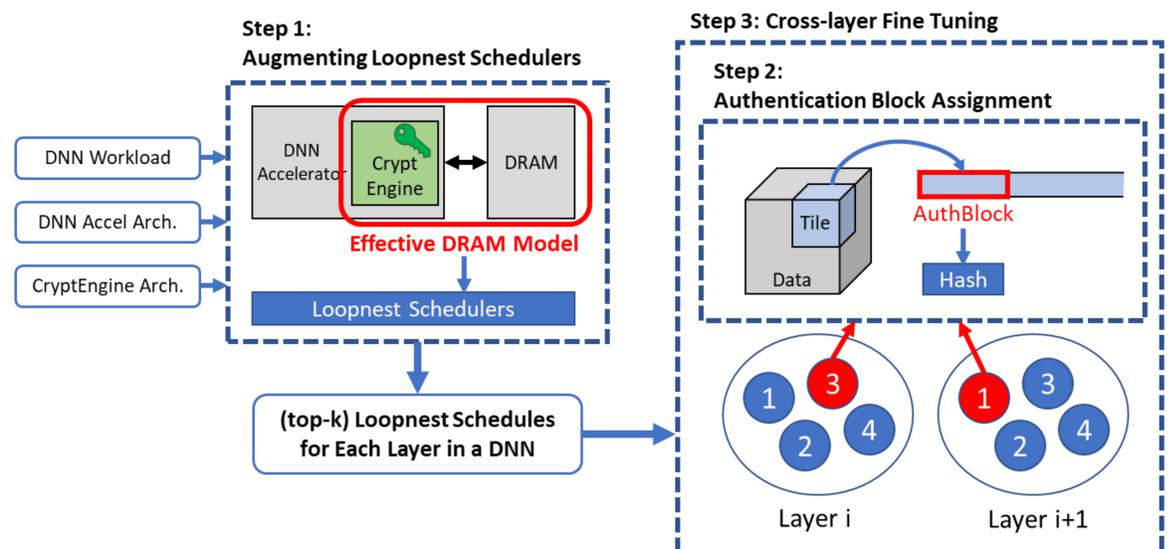
Run exhaustive search for possible sizes and orientations of AuthBlocks

Joint Search with Annealing

Define a neighbor as top-k loopnests for each layer
Visit random layer and neighbor at each iteration

- 1: $L_1, \dots, L_n \leftarrow L_1^o, \dots, L_n^o$
- 2: $cost \leftarrow PerfModel(L_1, \dots, L_n)$
- 3: $t \leftarrow T_{init}$ {initialize temperature}
- 4: **for** $n \leftarrow 1, \dots, N$ **do**
- 5: $i \leftarrow random(1, \dots, n)$
- 6: $L'_i \leftarrow GetNeighbor(L_i)$
- 7: $cost' \leftarrow PerfModel(L_1, \dots, L'_i, \dots, L_n)$
- 8: $cost_diff = cost - cost'$
- 9: **if** $exp \frac{cost_diff}{t} > random.uniform(0, 1)$ **then**
- 10: $L_i \leftarrow L'_i$ {probabilistic accept the new schedule}
- 11: $cost \leftarrow cost'$
- 12: **end if**
- 13: $t \leftarrow GetTemperature(t, n, T_{init}, T_{final})$
- 14: **end for**

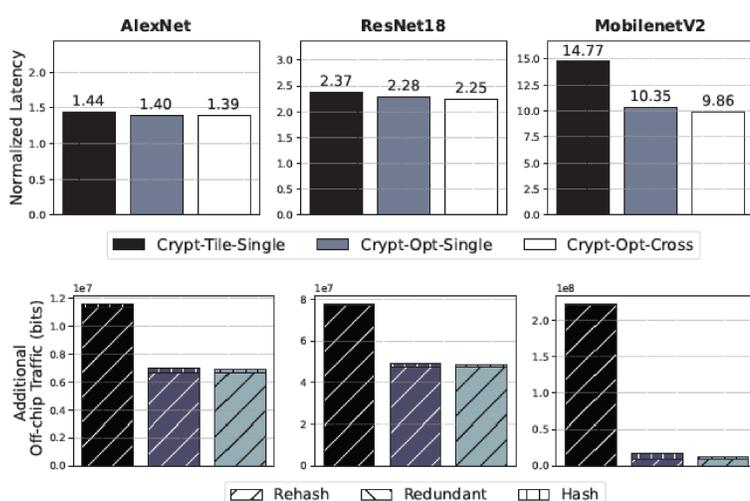
Putting All Together



Comparing Scheduling Algorithms

Comparing scheduling algorithms

- ~33% faster, ~50% better in EDP compared to the baseline
- Removes rehashing and identifies the optimal assignment
- Improvement larger for deeper workloads



Key Takeaways

Goal

Systematic investigation of the performance, area, and energy trade-off of secure DNN accelerators

SecureLoop Contributions

- Design space exploration framework for secure DNN accelerators supporting a TEE
- Scheduling for secure accelerators include authentication block assignment
- Cross-layer dependency considered from the loopnest schedule level
- Overall ~33% faster, ~50% better in EDP