

















PHD: Parallel Huffman Decoder on FPGA for Extreme Performance and Energy Efficiency

Yunkun Liao^{1,2,3}, Jingya Wu¹, Wenyan Lu^{1,4}, Xiaowei Li^{1,3}, Guihai Yan^{1,4}

SKLP, Institute of Computing Technology, Chinese Academy of Sciences¹

University of Chinese Academy of Sciences²

Zhongguancun Laboratory³

YUSUR Technology Co., Ltd.4

















CONTENTS



Background Knowledge



Experiment Results



Proposed Architecture: PHD



Summary





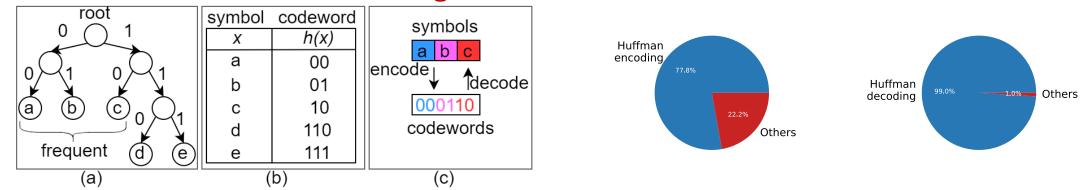
Background Knowledge





Huffman Coding

- A lossless data compression scheme.
 - Main idea: Encodes frequent symbols using shorter codewords and vice visa.
- Widely used in many data compression algorithms.
 - LZSS, ZIP, JPG, etc.
- Domain-specific accelerator design for Huffman coding is required.
 - This work accelerates Huffman decoding.



(a) Huffman coding tree; (b) Huffman coding table; (c) Huffman encoding Profiling cor and decoding.

Profiling conducted on the LZSS compression¹

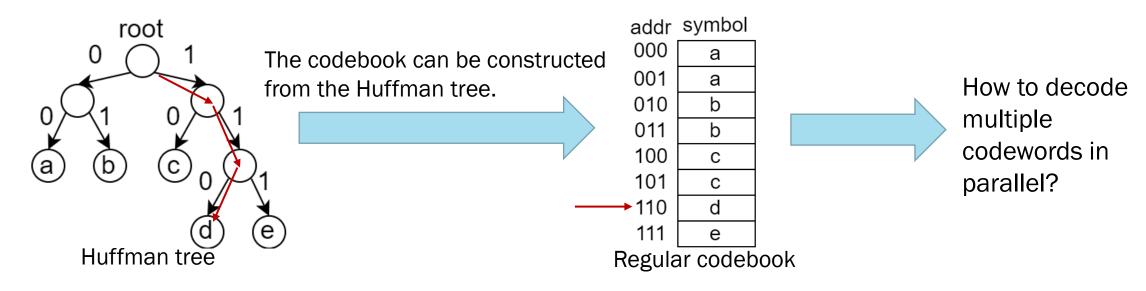
An accelerator for Huffman decoding is required.





Huffman Decoding

- Bit-serial decoding based on traversing the Huffman tree.
- Bit-parallel decoding based on looking up codebook.
 - Dominated method in existing Huffman decoding accelerators.



Decoding "110" to "d" requires three traversals.

Decoding "110" to "d" requires one lookup.

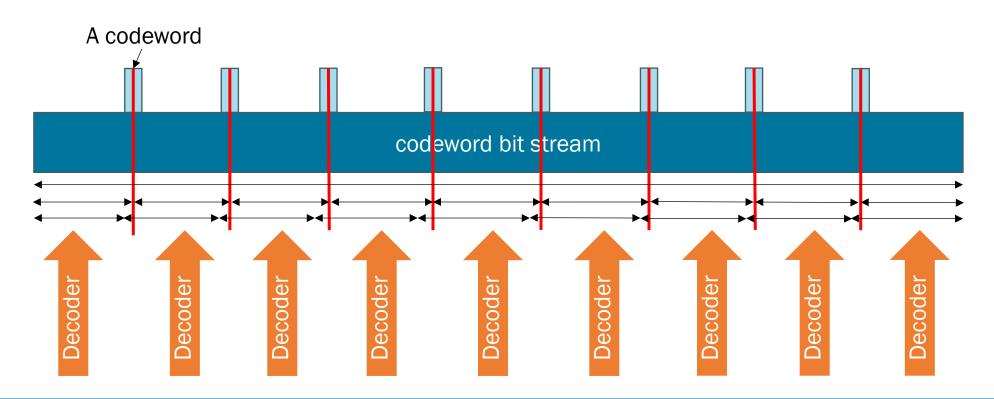
Huffman decoding requires codeword-level parallelism.





Parallelizing Decoding is Hard

Codeword-level parallelism can be obtained by partitioning the codeword bit stream.



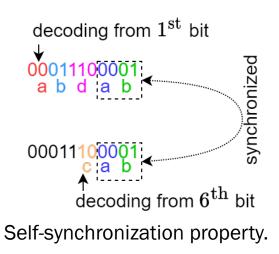
Parallelizing the decoding requires identifying codeword boundaries.

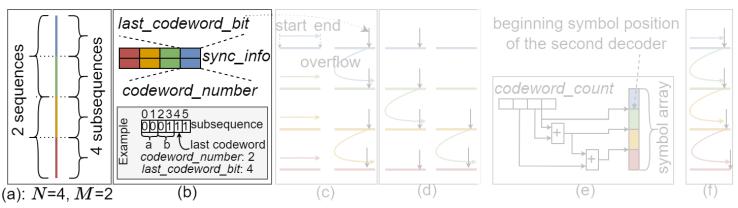




Parallelizing Decoding by Self-synchronization

- Almost all Huffman codes have the self-synchronization property².
- Achieve codeword-level parallelism with the self-synchronization.
 - Phase_1: Intra-sequence synchronization.
 - Phase_2: Inter-sequence synchronization.
 - Phase_3: Calculating the codeword boundaries.
 - Phase_4: Parallel decoding on N subsequences.





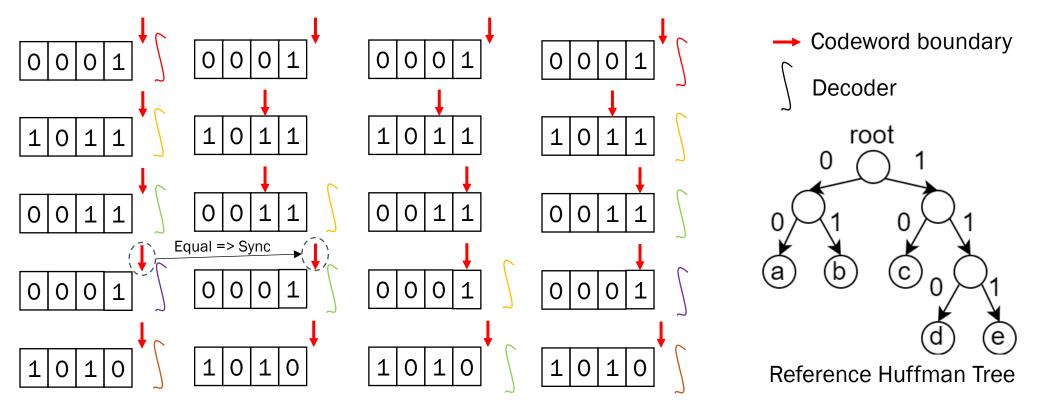
(a) Codewords splitting; (b) Synchronization information; (c) Phase_1; (d) Phase_2 (e) Phase_3 (f) Phase_4.

Self-synchronization trades computational complexity for parallelism.





Parallelizing Decoding Example



Decoder state table: Sync $\sqrt{\ }$, Not Sync \times \circ

S	×	V	√
3	×	×	√
3	×	×	\checkmark
2	×	√	V
7	×	V	V

Assume decoding one symbol requires constant time T

Serial decoding: 9T

Parallel decoding: 8T





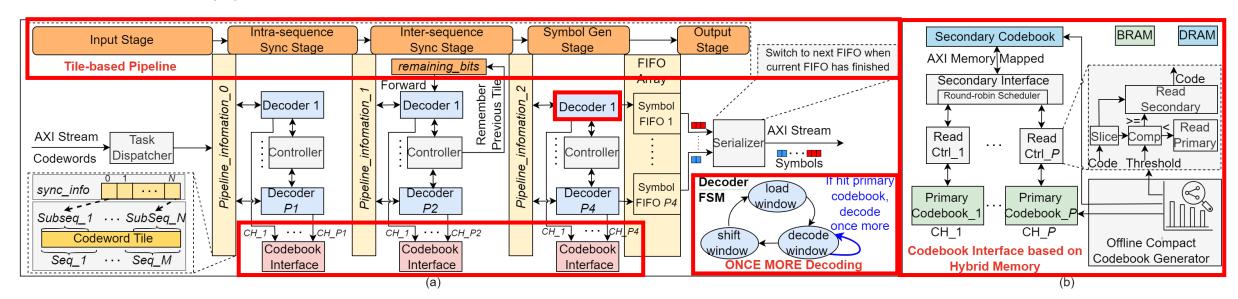
Proposed Architecture: PHD





Overall Architecture of PHD

- PHD: The first self-synchronization-based Parallel Huffman Decoder architecture for FPGA.
 - Codebook interface based on hybrid memory.
 - ONCE MORE decoding optimization.
 - Tile-based pipeline dataflow.



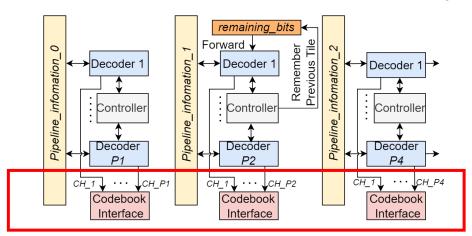
PHD is the first accelerator to leverage the self-synchronization.





The Challenge of Codebook Storage

- The decoders of PHD is based on bit-parallel Huffman decoding, requiring codebook.
- Homogeneous storage designs.
 - BRAM-only: Fast, but with limited codebook size.
 - Regular codebook for S-bit Huffman code is 2^S bytes.
 - For HTTP/2 HPACK Huffman code³, S is 30, requiring 2GB.
 - DRAM-only: Support large codebook, but slow.
- Low-latency and storage-efficient codebook interface is required.



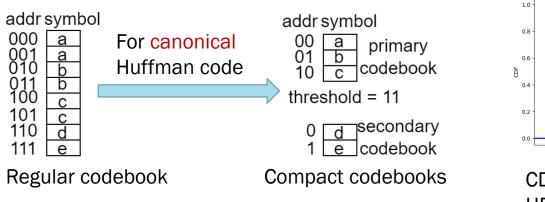
PHD should balance the lookup speed and BRAM consumption.

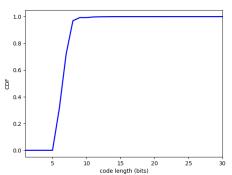




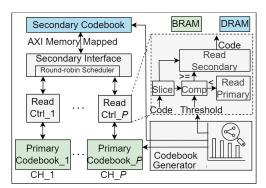
Codebook Interface based on Hybrid Memory

- Key enabler: Most Huffman codes are canonical. The regular codebook for canonical code can be split into one primary codebook and multiple secondary codebooks⁴.
- Our insight: Primary codebook is for short codewords. The symbols of short codeword are more common in Huffman code by design.
- Our hybrid storage design:
 - Multi-channel BRAM for primary codebook
 - Shared-channel DRAM for secondary codebooks.





CDF of code length from HPACK-Huffman files



Codebook interface based on hybrid memory

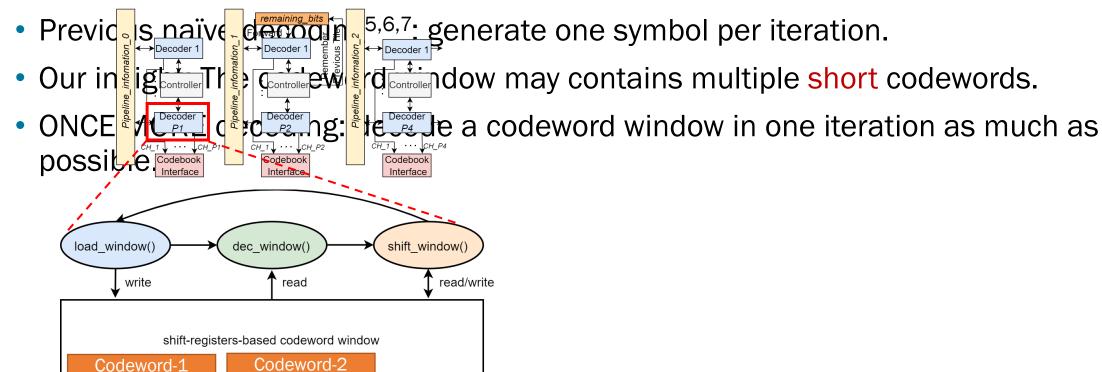
PHD makes short (common) -codeword lookup fast.





ONCE MORE Decoding

Codewords are sent to a shift-register-based window for decoding.



The shift register can contain the longest codeword.

PHD makes subsequence-level decoders fast.





The Challenge of Codeword-Tile Dependency

- On-chip codeword tile is used to support one-cycle codewords access.
- The size of on-chip codeword tile is limited and codewords are split into multiple tiles.
- There are dependencies between consecutive tiles.
 - A codeword may span two tiles.
- Ignoring the dependencies between consecutive tiles leads to incorrect decoding results
 of the remaining codeword tiles.

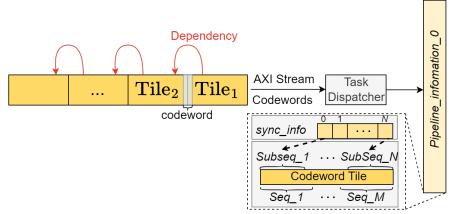


Illustration of the codeword-tile dependency

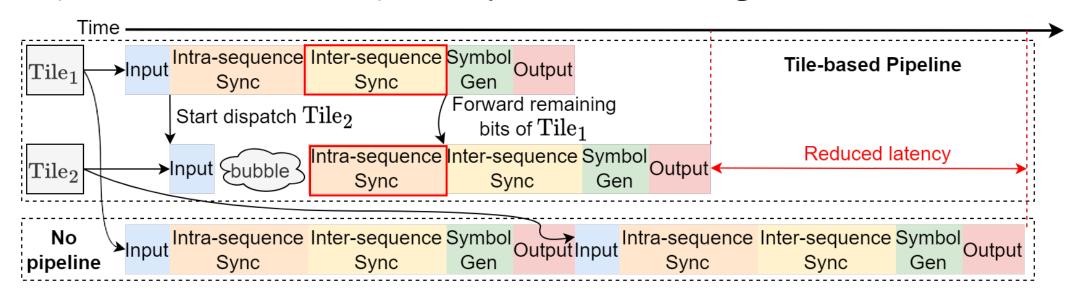
PHD should guarantee the correctness.





Tile-based Pipeline

- Our insight: no dependencies between codeword tiles at the subsequence level in the intra-sequence synchronization stage.
 - Intra-sequence Sync of Tile₂ can run in parallel with the inter-sequence Sync of Tile₁.
- Our methods: remember and forward the remaining bits of the previous tile to the subsequent tile in the inter-sequence synchronization stage.



PHD enables tile-level parallelism without compromising correctness.





Experiment Results





Setup

- Huffman code: HTTP/2 HPACK, used for compressing the HTTP/2 packet header fields.
- PHD: Described using the Xilinx HLS C++, Xilinx Vitis 2022.1 is used for cycle-accurate simulations, synthesis and implementations.
- Baselines:
 - CPU: a performance-optimized library called LS-HPACK⁸, running on Intel Core-i5 12400 processor.
 - GPU: the open-source code of Weißenberger⁹, adapted to HTTP/2 HPACK, running on NVIDIA GeForce RTX 4090 GPU.
 - FPGA: Bit-parallel, described using the Xilinx HLS C++.

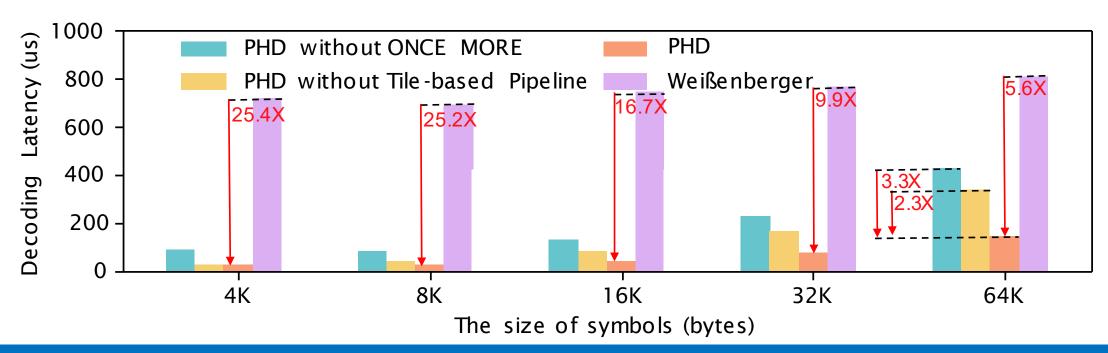




Decoding Latency

Results:

- Compared with CPU baseline: a latency reductions ranging from 24.5X to 60.6X.
- Compared with bit-parallel FPGA baseline: a latency reduction ranging from 18.5X to 57.5X.
- Compared with GPU baseline: latency reductions of 5.6X to 25.4X over Weißenberger.



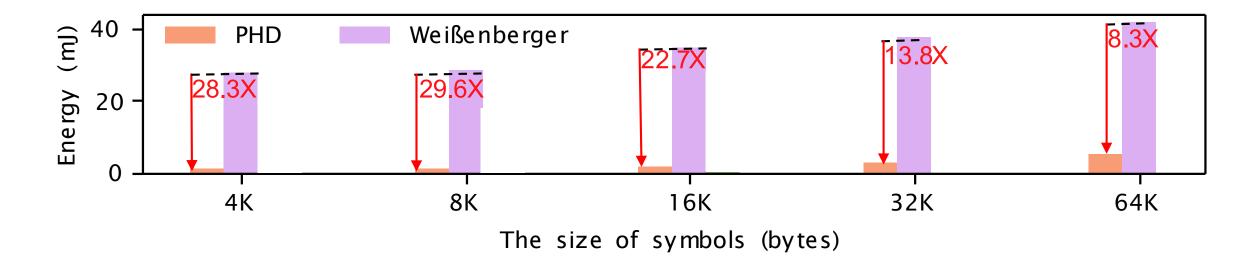
PHD is faster!





Energy Consumption

- Results:
 - PHD reduces energy by 8.3X to 29.6X over Weißenberger.



PHD is more energy-efficient!



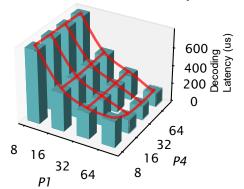


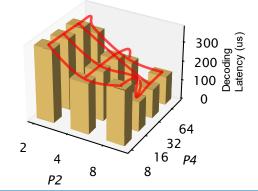
Design Space Exploration

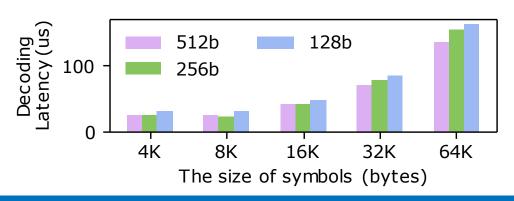
- Architecture parameters in PHD:
 - P1: Subsequence-level parallelism in intra-sequence synchronization stage.
 - P2: Subsequence-level parallelism in the inter-sequence synchronization stage.
 - P4: Subsequence-level parallelism in the symbol generation stage.
 - Subsequence size: workload capacity of each decoder.

Results:

- Large P1 and P2 parallelism are crucial for overall performance.
- Smaller subsequence size leads to a little longer latency







PHD is parameterized!



Summary





Summary

- We presented PHD, the first accelerator targeting at self-synchronizationbased parallel Huffman decoding.
 - PHD realizes bit-level, codeword-level and tile-level parallelism.
 - PHD implements a compact codebook interface based on hybrid memory.
 - PHD proposes ONCE MORE optimization to accelerate subsequence decoding.
- Future work:
 - An analytical performance model of PHD.
 - Generating and tuning PHD automatically.

If you need the code for research, please send an email to me! ©





SHAPING THE NEXT GENERATION OF ELECTRONICS

JUNE 23-27, 2024

MOSCONE WEST CENTER SAN FRANCISCO, CA, USA

THANK YOU!













