#### Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

#### Amirali Boroumand

Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, Onur Mutlu









SEOUL NATIONAL UNIVERSITY





#### **Consumer Devices**



#### **Consumer devices are everywhere!**

#### Energy consumption is a first-class concern in consumer devices





#### Popular Google Consumer Workloads



Chrome

**Google's web browser** 



#### **TensorFlow Mobile**

Google's machine learning framework



**Google's video codec** 



# **Energy Cost of Data Movement**

#### I<sup>st</sup> key observation: 62.7% of the total system energy is spent on data movement



**Processing-in-Memory (PIM)** 

Potential solution: move computation close to data

Challenge: limited area and energy budget

#### Using PIM to Reduce Data Movement

2<sup>nd</sup> key observation: a significant fraction of data movement often comes from simple functions

We can design lightweight logic to implement these <u>simple functions</u> in <u>memory</u>

Small embedded low-power core



Small fixed-function accelerators



Offloading to PIM logic reduces energy by 55.4% and improves performance by 54.2% on average

#### Goals

Understand the data movement related bottlenecks in modern consumer workloads

2 Analyze opportunities to mitigate data movement by using processing-in-memory (PIM)

**3** Design PIM logic that can maximize energy efficiency given the limited area and energy budget in consumer devices

# Outline

- Introduction
- Background
- Analysis Methodology
- Workload Analysis
- Evaluation
- Conclusion

#### **Potential Solution to Address Data Movement**

- Processing-in-Memory (PIM)
  - A potential solution to reduce data movement
  - Idea: move computation close to data
    - **Reduces data movement**
    - **Exploits large in-memory bandwidth**
    - **V** Exploits shorter access latency to memory
- Enabled by recent advances in 3D-stacked memory



# Outline

- Introduction
- Background
- Analysis Methodology
- Workload Analysis
- Evaluation
- Conclusion

# Workload Analysis Methodology

- Workload Characterization
  - Chromebook with an Intel Celeron SoC and 2GB of DRAM
  - Extensively use performance counters within SoC
- Energy Model
  - Sum of the energy consumption within the CPU, all caches, off-chip interconnects, and DRAM



VP9

# **PIM Logic Implementation**



#### SAFARI

## Workload Analysis



Chrome

**Google's web browser** 



**TensorFlow** 

Google's machine learning framework





## Workload Analysis



#### Chrome

**Google's web browser** 



VP9 VouTube Video Playback Google's video codec



#### How Chrome Renders a Web Page



#### SAFARI

#### How Chrome Renders a Web Page



## **Browser Analysis**

- To satisfy user experience, the browser must provide:
  - Fast loading of webpages
  - Smooth scrolling of webpages
  - Quick switching between browser tabs
- We focus on two important user interactions:
  - I) Page Scrolling
  - 2) Tab Switching
  - Both include page loading

## Scrolling



#### What Does Happen During Scrolling?



# Scrolling Energy Analysis



41.9% of page scrolling energy is spent on texture tiling and color blitting

#### SAFARI

# Scrolling a Google Docs Web Page



## Scrolling a Google Docs Web Page



Can we use PIM to mitigate the data movement cost for texture tiling and color blitting?





Texture tiling is a good candidate for PIM execution





#### Can We Implement Texture Tiling in PIM Logic?



Requires simple primitives: memcopy, bitwise operations, and simple arithmetic operations



9.4% of the area available for PIM logic

7.1% of the area available for PIM logic

PIM

**Accelerator** 

PIM core and PIM accelerator are feasible to implement in-memory Texture Tiling

## **Color Blitting Analysis**

Generates a large amount of data movement Accounts for 19.1% of the total system energy during scrolling

#### Color blitting is a good candidate for PIM execution

Requires low-cost operations: Memset, simple arithmetic, and shift operations

It is feasible to implement color blitting in PIM core and PIM accelerator



# Scrolling Wrap Up

Texture tiling and color blitting account for a significant portion (41.9%) of energy consumption

**37.7%** of total system energy goes to data movement generated by these functions

Both functions can benefit significantly from PIM execution

Both functions are feasible to implement as PIM logic



#### **Tab Switching**



#### What Happens During Tab Switching?

- Chrome employs a multi-process architecture
  - Each tab is a separate process



- Main operations during tab switching:
  - Context switch
  - Load the new page

#### SAFARI

# **Memory Consumption**

- Primary concerns during tab switching:
  - How fast a new tab loads and becomes interactive
  - Memory consumption

# Chrome uses compression to reduce each tab's memory footprint



### Data Movement Study

• To study data movement during tab switching, we emulate a user switching through 50 tabs

We make two key observations:

Compression and decompression contribute to 8.8% of the total system energy

> **I 9.6 GB** of data moves between CPU and ZRAM

2

#### Can We Use PIM to Mitigate the Cost?



PIM core and PIM accelerator are feasible to implement in-memory compression/decompression

#### SAFARI

# Tab Switching Wrap Up

A large amount of data movement happens during tab switching as Chrome attempts to compress and decompress tabs

Both functions can benefit from PIM execution and can be implemented as PIM logic

## Workload Analysis



Chrome

**Google's web browser** 



**TensorFlow** 

Google's machine learning framework





## Workload Analysis



Chrome Google's web browser



Google's machine learning framework



SAFARI



# TensorFlow Mobile





## Packing



A simple data reorganization process that requires simple arithmetic

#### Quantization



A simple data conversion operation that requires shift, addition, and multiplication operations

#### SAFARI





A simple data conversion operation that requires shift, addition, and multiplication operations

#### SAFARI

## Video Playback and Capture





#### Majority of energy is spent on data movement

# Majority of data movement comes from simple functions in decoding and encoding pipelines



# Outline

- Introduction
- Background
- Analysis Methodology
- Workload Analysis
- Evaluation
- Conclusion

## **Evaluation Methodology**

- System Configuration (gem5 Simulator)
  - SoC: 4 OoO cores, 8-wide issue, 64 kB Ll cache, 2MB L2 cache
  - **PIM Core:** I core per vault, I-wide issue, 4-wide SIMD, 32kB LI cache
  - **3D-Stacked Memory: 2GB cube, 16 vaults per cube** 
    - Internal Bandwidth: 256GB/S
    - Off-Chip Channel Bandwidth: 32 GB/s
  - Baseline Memory: LPDDR3, 2GB, FR-FCFS scheduler
- We study each target in isolation and emulate each separately and run them in our simulator **40** SAFARI

## **Normalized Energy**



and packing comes from eliminating data movement

PIM core and PIM accelerator reduces energy consumption on average by 49.1% and 55.4% SAFARI

#### **Normalized Runtime**

#### CPU-Only PIM-Core PIM-Acc



Offloading these kernels to PIM core and PIM accelerator improves performance on average by 44.6% and 54.2%

## Conclusion

- Energy consumption is a major challenge in consumer devices
- We conduct an in-depth analysis of popular Google consumer workloads
  - 62.7% of the total system energy is spent on data movement
  - Most of the data movement often come from <u>simple functions</u> that consists of <u>simple operations</u>
- We use **PIM** to reduce data movement cost
  - We design lightweight logic to implement simple operations in DRAM



- Reduces total energy by 55.4% on average
- Reduces execution time by 54.2% on average

#### SAFARI

#### Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

#### Amirali Boroumand

Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, Onur Mutlu









SEOUL NATIONAL UNIVERSITY









## Video Playback



63.5% of the system energy is spent on <u>data movement</u>

80.4% of the data movement energy comes from <u>sub-pixel interpolation</u> and <u>deblocking filter</u>

Sub-pixel interpolation: interpolates the value of pixels at non-integer location SAFARI

Deblocking filter: a simple lowpass filter that attempts to remove discontinuity in pixels 39

#### Video Playback



#### Based on our analysis, we conclude that:

- Both Functions can benefit from PIM execution
- They can be implemented in PIM logic

sup-pixer inter polation and deblocking inter

<u>Sub-pixel interpolation</u>: interpolates the value of pixels at non-integer location SAFARI

Deblocking filter: a simple lowpass filter that attempts to remove discontinuity in pixels

# Video Capture

SAFARI



59.1% of the system energy is spent on data movement

Majority of the data movement energy comes from motion estimation which accounts for 21.3% of total system energy

<u>Motion estimation</u>: compresses the frames using temporal redundancy between them

## Video Capture



#### Motion Estimation is a good candidate for PIM execution and is feasible to be implemented at PIM logic

Majority of the data movement energy comes from motion estimation which accounts for 21.3%

of total system energy

**Motion estimation**: compresses the frames using temporal redundancy between them



# TensorFlow Mobile

57.3% of the inference energy is spent on <u>data movement</u>

54.4% of the data movement energy comes from packing/unpacking\_and quantization

Packing: A simpleQuantization: A simple datadata re-organization processconversion operation thatthat reorders elements of matrice32-bit floating point to 8-SAFARIbit integers35

#### **TensorFlow Mobile**



#### Based on our analysis, we conclude that:

- Both Functions are good candidate for PIM execution
- It is feasible to implement them in PIM logic

54.4% of the data movement energy comes from packing/unpacking\_and quantization

Packing: A simpleQuantization: A simple datadata re-organization processconversion operation that convertsthat reorders elements of matrices32-bit floating point to 8-bitSAFARIintegers32

## Data Movement Study

- To study data movement during tab switching, we did an experiment:
  - A user opens 50 tabs (most-accessed websites)
  - <u>Scrolls</u> through <u>each</u> for a few second
  - <u>Switches</u> to the <u>next tab</u>
- We make two key observations:
  - Compression and decompression contribute to 18.1% of the total system energy
  - 19.6 GB of data swapped in and out of ZRAM

#### **Other Details and Results in the paper**

#### Detailed discussion of the other workloads

- Tensorflow Mobile
- Video Playback
- Video Capture
- Detailed discussion of VP9 hardware decoder and encoder
- System Integration
  - Software interface
  - Coherence

### **Color Blitting Analysis**

- Color blitting generates a large amount of data movement
  - Accounts for 19.1% of the total system energy used during page scrolling
  - 63.9% of the energy consumed by <u>color blitting</u> is due to data movement

#### **Based on our analysis, we conclude that:**

- <u>Color blitting</u> is <u>a good candidate</u> for <u>PIM execution</u>
- It is feasible to implement color blitting in PIM logic
- Color blitting require low-cost computation operations
  - <u>Memset</u>, <u>simple arithmetic</u> for alpha blending, <u>shift</u> operations



#### **Data Movement During Tab-Switching**



#### Methodology: Implementing PIM Targets



Gener Eirfedrussetip MPLM Accelerator

- A custom zest on zest of the single store and a budget for the A 4-wide storm subject single thread of the Mogiet

#### SAFARI

## **Cost Analysis: Texture Tiling**



# **Memory Consumption**

#### Potential Solution: kill <u>inactive background</u> tabs

- When the user <u>accesses those tabs</u> again, reload the tab pages from the disk
- Downsides: I) high latency 2) loss of data



Aw, Snap!

#### Chrome uses compression to reduce <u>each tab's</u> <u>memory footprint</u>

- Uses <u>compression</u> to compress <u>inactive pages</u> and places them into a DRAM-based memory pool (ZRAM)
- When the user switches to a previously-inactive tab, Chrome loads the data from ZRAM and decompresses it

#### SAFARI

# Methodology: Identifying PIM targets

- We pick a function as <u>PIM target</u> candidate if:
  - Consumes the most energy out of the all functions
  - Significant <u>data movement</u> cost (MPKI > 10)
  - Bounded by data movement, not computation
- We drop any candidate if:
  - X Incurs <u>any performance loss</u> when runs on <u>PIM logic</u>
  - **X** Requires more area than is available in the logic layer of 3Dstacked memory