ApproxEye: enabling approximate computation reuse for micro-robotic computer vision

> Xin He, Guihai Yan, Faqiang Sun, Yinhe Han and Xiaowei LI Speaker: Yinhe Han

State Key Laboratory of Computer Architecture, Institute of Computing Technology , Chinese Academy of Sciences





Inefficiency of nowadays micro-sized system

- Existing micro-sized computing system (e.g. micro-robotics) fails to meet targeted performance requirement
 - Due to:



Energy efficient computing



Prior ways to boost system performance

- To perform complex tasks like rescuing, tracking abnormal object tracking and so on
 - Computation Offloading
 - Responsiveness and bandwidth
 - Customized accelerator
 - Portability issue
 - Energy efficient scheduling
 - Limited achievable speedup







Approximate Computing(Domain specific approach)

Leveraging intrinsic application resilience to improve efficiency

E.g. DCT Image Compression



Original method Approximate method

By using inexact multiplier

Outcome: **2.82%** reduction on PSNR Performance: **31.44%** Improvement





By reducing excess iterations

Outcome: **1.3%** accuracy reduction Performance: **4.97X** Improvement

Approximate Computing opens another road towards computing efficiency

Targeting and observation

Focusing on accelerating the computer vision algorithm through approximate computing:



We select optical flow to conduct approximate



Optical flow as target killer application

- We select optical flow to conduct approximate computing
 - Optical flow is <u>widely used</u> for visual surveillance, motion estimation, object tracking.
 - Its computation kernel is a <u>typical</u> kind of operations seen in computer vision



> To reuse the prior computations:

- Exact reuse scheme (100% accurate)
 - Traditional computation reuse
 - Fail to reuse similar computations
- Approximate reuse scheme (allow accuracy degradation)
 - RACB scheme
 - Fail to fully extract reuse opportunities
 - Proposed ApproxEye scheme
 - Significance aware computing
 - Adaptive reuse granularity selection

For traditional reuse scheme

Traditional computation reuse scheme

- Leveraging computation locality to reuse computations
 - But the tight requirement is used for 100% precision



For traditional reuse scheme

Traditional computation reuse scheme

- Leveraging computation locality to reuse computations
 - Require all inputs of a reused item is exactly equal to underlying computation





Limitation of prior approximate reuse scheme

Prior approx reuse (e.g. RACB scheme [islped2005])

• Reuse requirement: the equality of <u>Most Significant Bits</u>



Cons :

- It's non-trivial to determine the length of LSBs mask
- A unified mask for all different inputs limits the effectiveness
- Limited application scope
- Predetermined reused granularity



- To fully extract the potential of computation reuse, ApproxEye proposes both
 - "Simple but effective" significance aware quantification
 - Adaptive reuse granularity selection
- Workflow of ApproxEye
 - Fully extract computation with high reuse possibilities
 - ② Based on statistical analysis, the optimal reuse granularity is specified
 - 3 Calculate input significance and set masks adaptively

- Fully extract computation with high reuse possibilities
 - Prohibit missing any reuse opportunities from history

Given: different optical flow kernel "multiply add" sequences, e.g. S1 and S1 Determine: similar substring inside both S1 and S2



- Once the similar computation is obtained, statistically analysis is conducted to
 - determine reuse granularity that gives highest speedup



Calculate input significance and set masks adaptively



This simple but faithful significance estimation benefits from: 1)This simple form can improve the effectiveness of runtime estimation 2)The significance can be easily used for adaptive LSB masking

创新成果二:可变粒度的近似计算方法

Parallel search tailored for approximate computing



Experiment

Experimental setup

- Target applications
 - Algorithm: State of the art Lucas Kanade pyramid iterative optical flow algorithm
 - Dataset: open source data from Microsoft redmond
- Performance and power simulation
 - Latency of multiply-add operations is obtained using Cadence tool flow in TMSC 45nm PTM model
 - The latency and power of TCAM is simulated from TCAM model based on Cacti

Experimental result

- The amount computation which can be reused by ApproxEye
 - Compared with oracle solution and RACB scheme



Experimental result

Speedup achieved by ApproxEye



Prior RACB approximate computing scheme :

- Unfaithful significance quantification using unified deviation thresholds for all different inputs
- Select the reuse granularity arbitrarily, miss the opportunities for partial reuse

Experimental result





Comparable or less power than Cache

Conclusion

- Approximate computing can exploit the potential of computation reuse in computer vision applications
- To fully extract the potential of computation reuse, one should take a deep look into data distribution and make comprehensive decisions
- Moreover, the engine of approximate computing should also be fine tuned for efficiency

Thanks for listening!

Questions?

Xin He

State Key Laboratory of Computer Architecture, Institute of Computing Technology, Chinese Academy of Sciences (ICT, CAS)



