

# Load Value Approximation

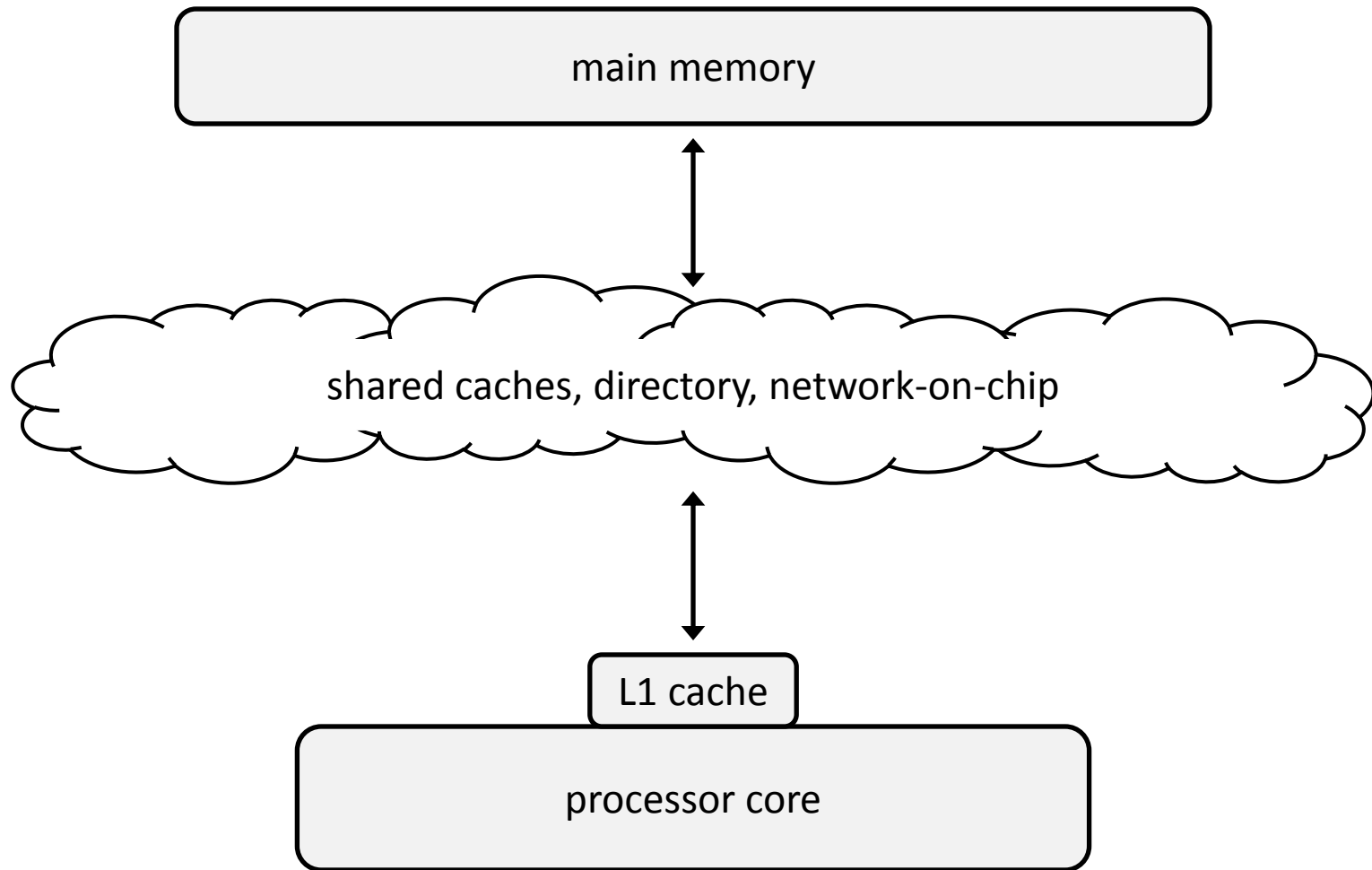
Joshua San Miguel

Mario Badr

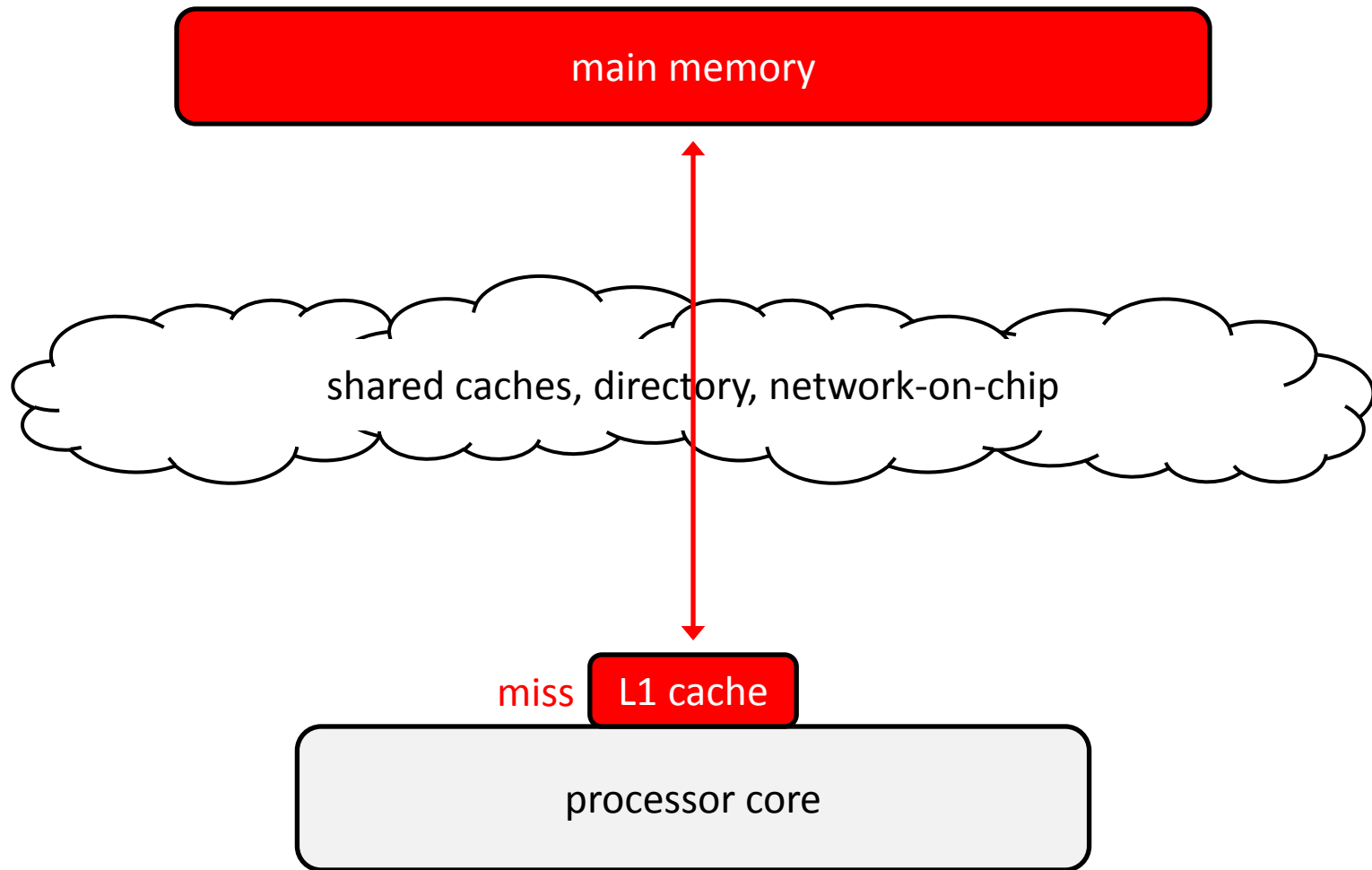
Natalie Enright Jerger



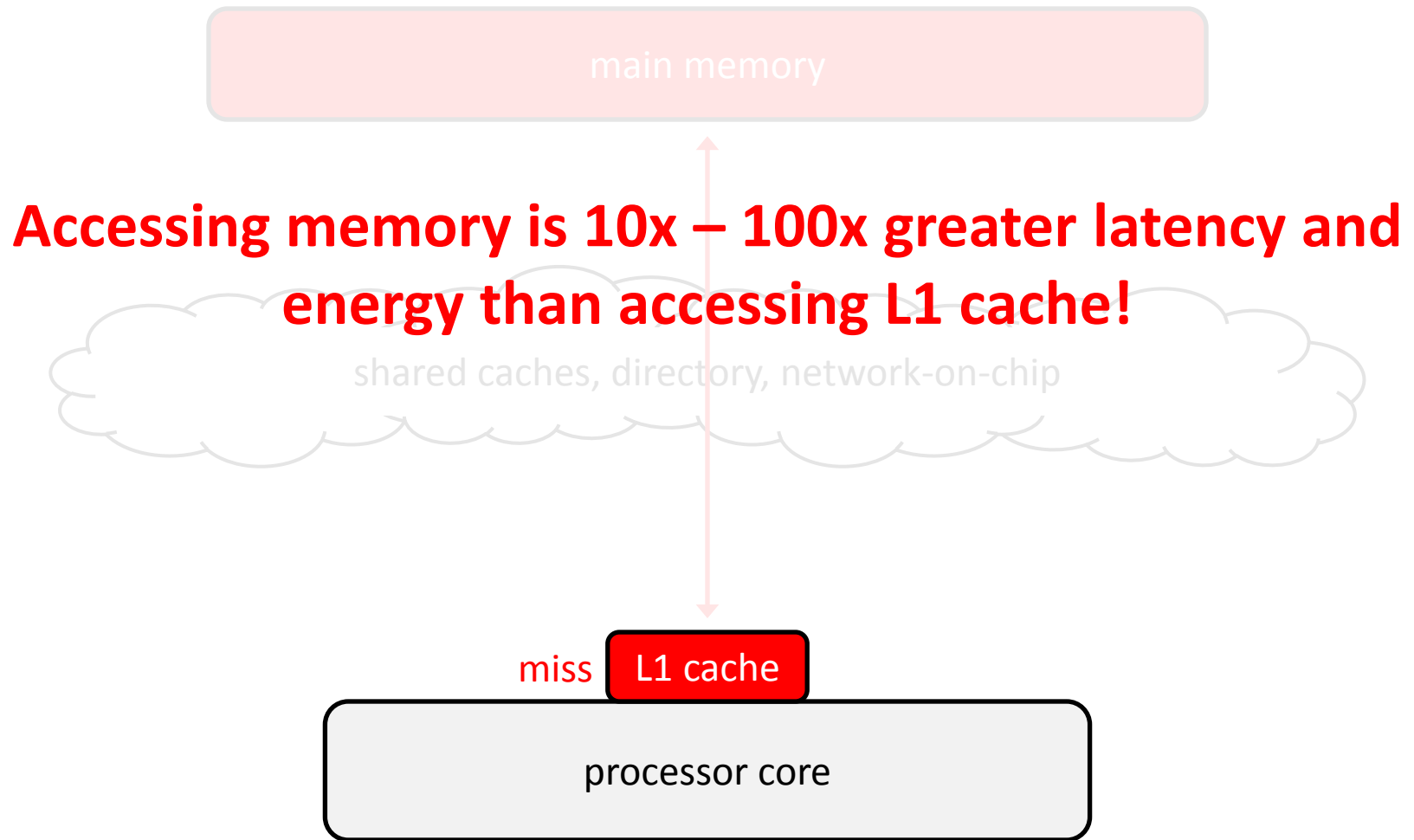
# Accessing Memory



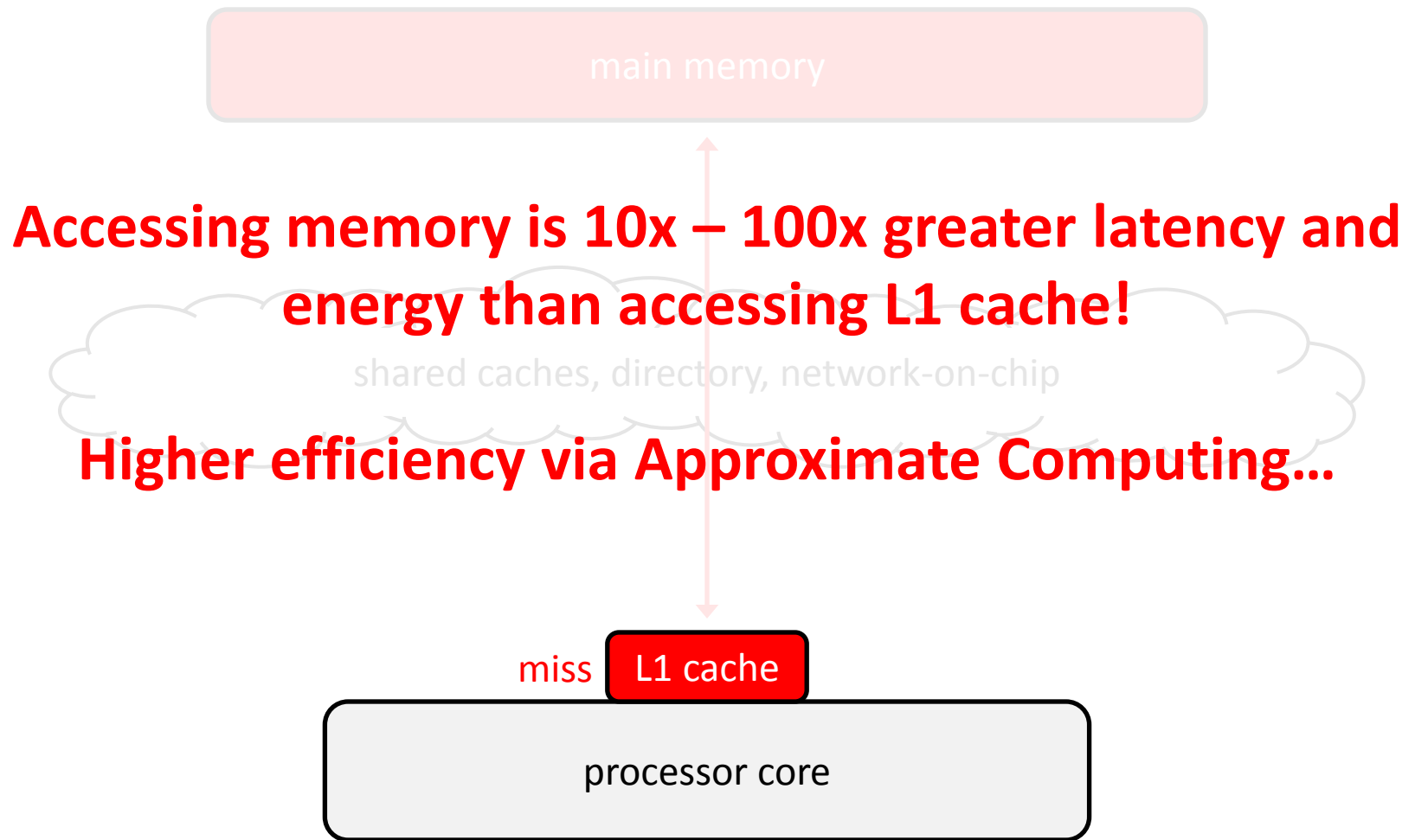
# Accessing Memory



# Accessing Memory



# Accessing Memory



# Approximate Computing

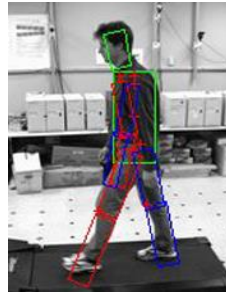
Not all computations need to be precise.

Data mining



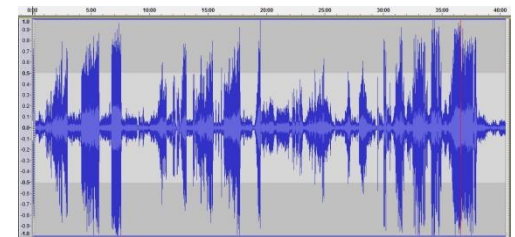
<http://www.zentut.com/>

Computer vision



<http://www.cc.gatech.edu/~cnieto6/>

Audio and video processing



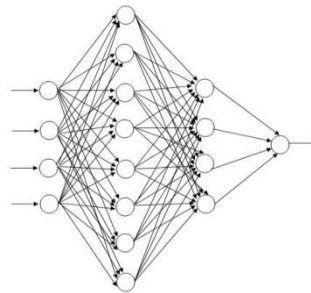
<http://themicparlour.blogspot.ca/>

Gaming



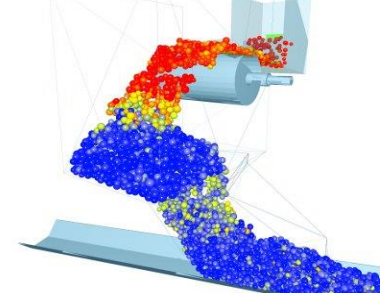
<http://www.businessweek.com/>

Machine learning



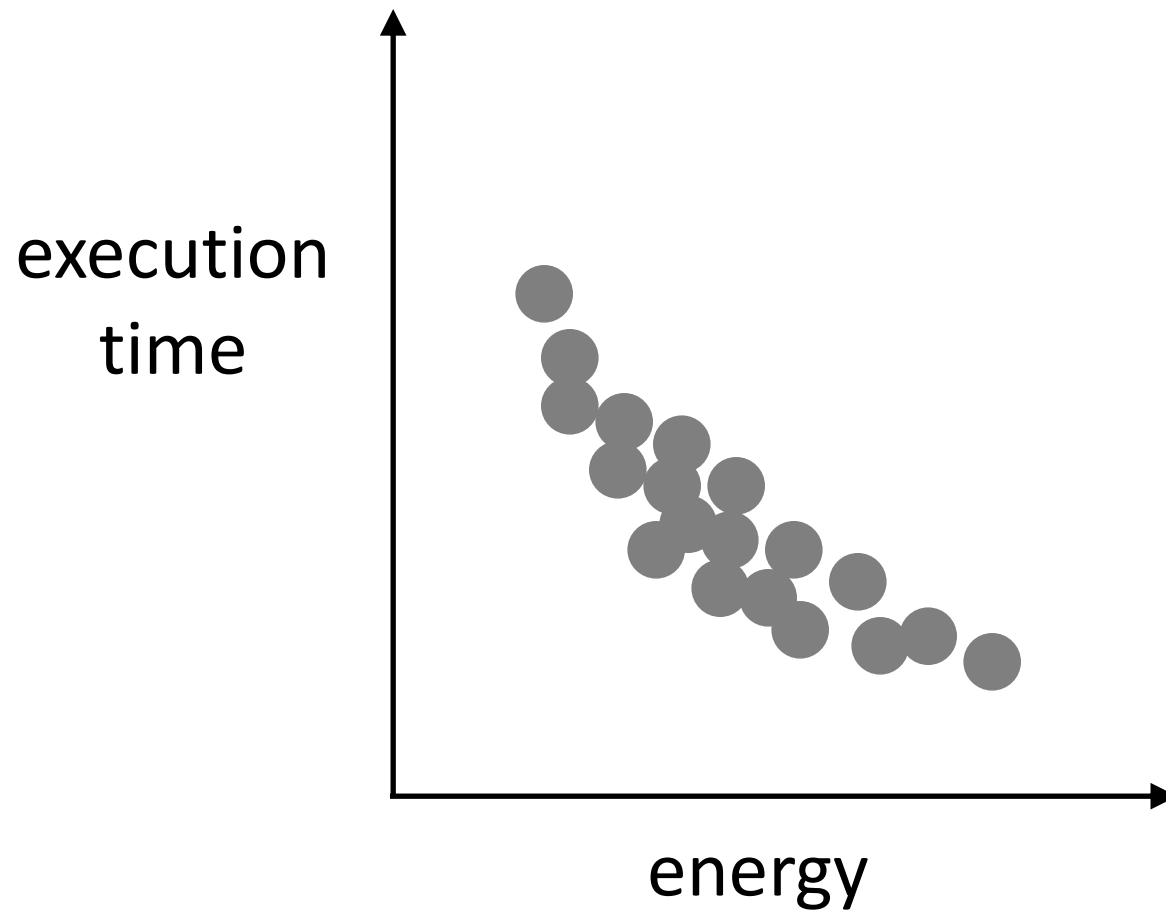
<http://www.analyticbridge.com/>

Dynamical simulation

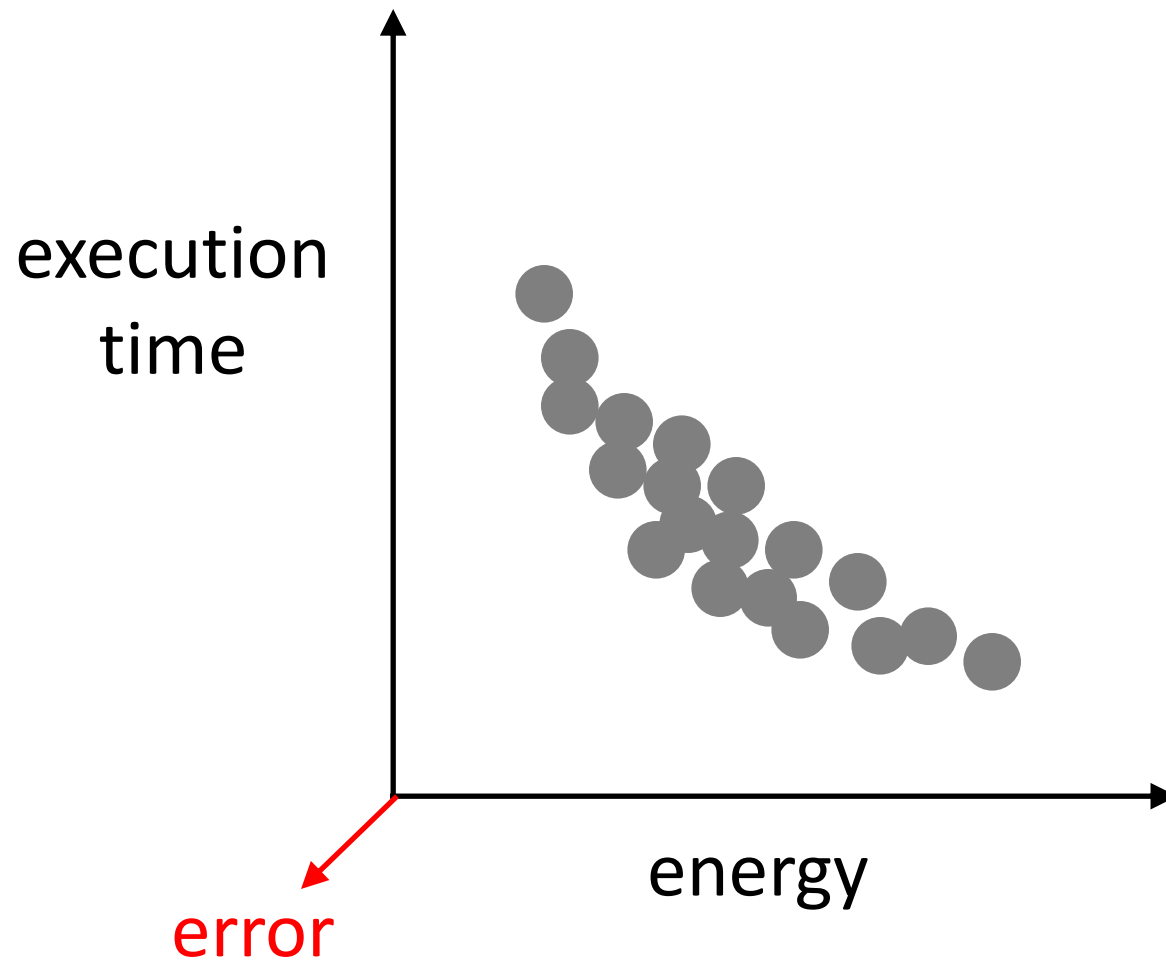


<http://www.scientific-computing.com/>

# Approximate Computing

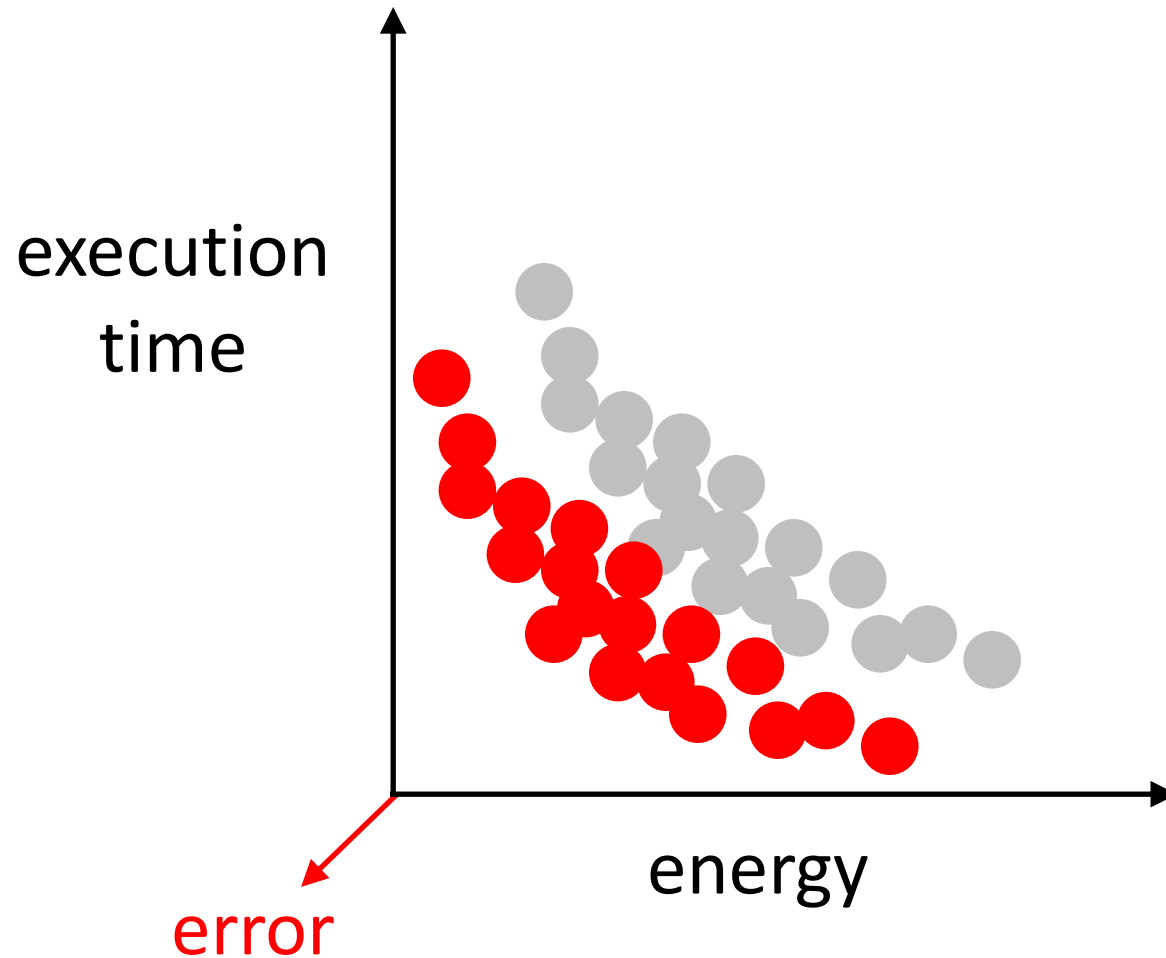


# Approximate Computing





# Approximate Computing



# Approximate Computing

Many applications can tolerate approximate data.

- 40% to nearly 100% of data footprint is approximate [Sampson, MICRO 2013].

# Approximate Computing

Many applications can tolerate approximate data.

- 40% to nearly 100% of data footprint is approximate [Sampson, MICRO 2013].

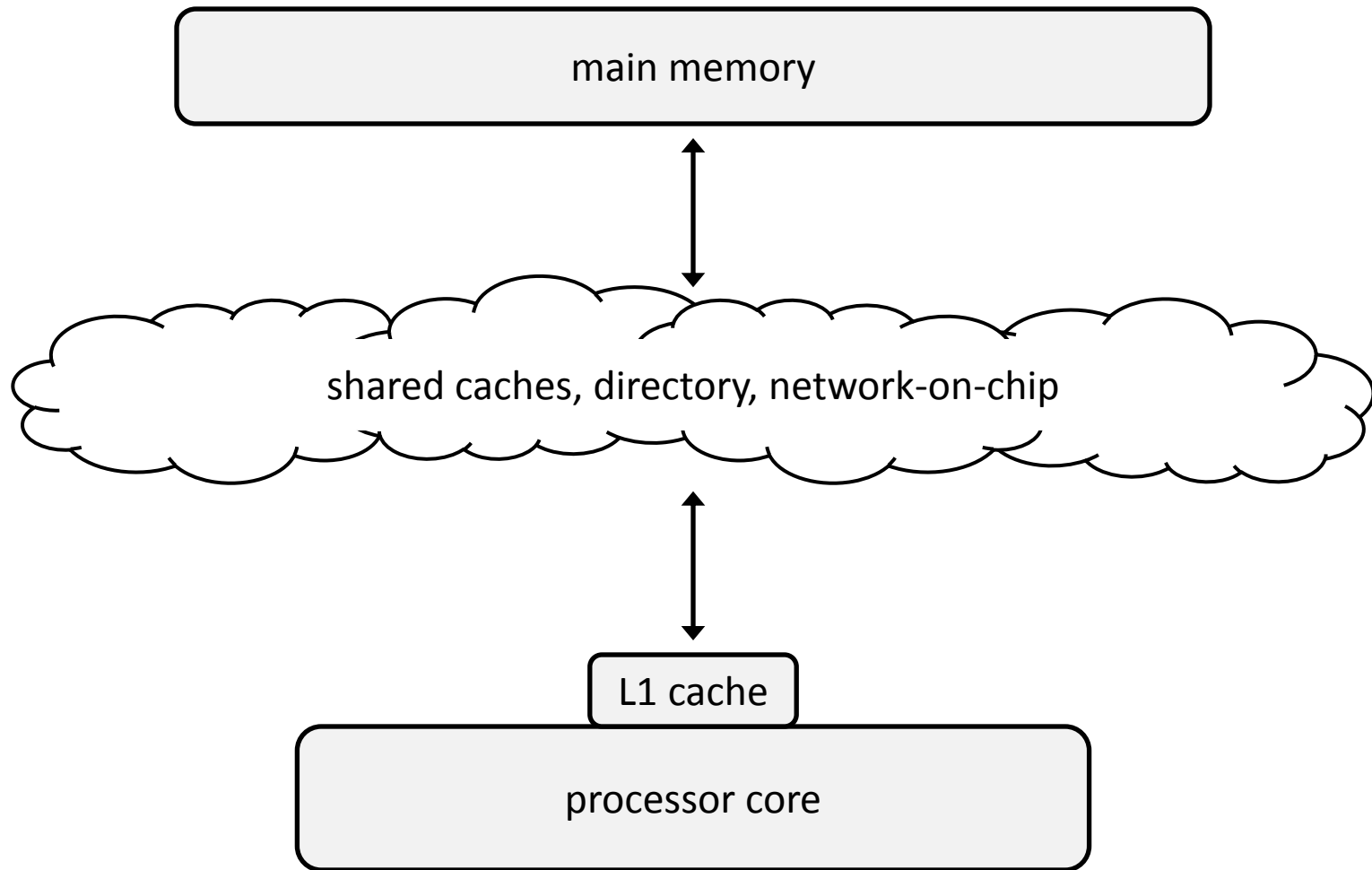
## Approximate value locality:

- Many data values are similar to or can be approximated from previously seen values.

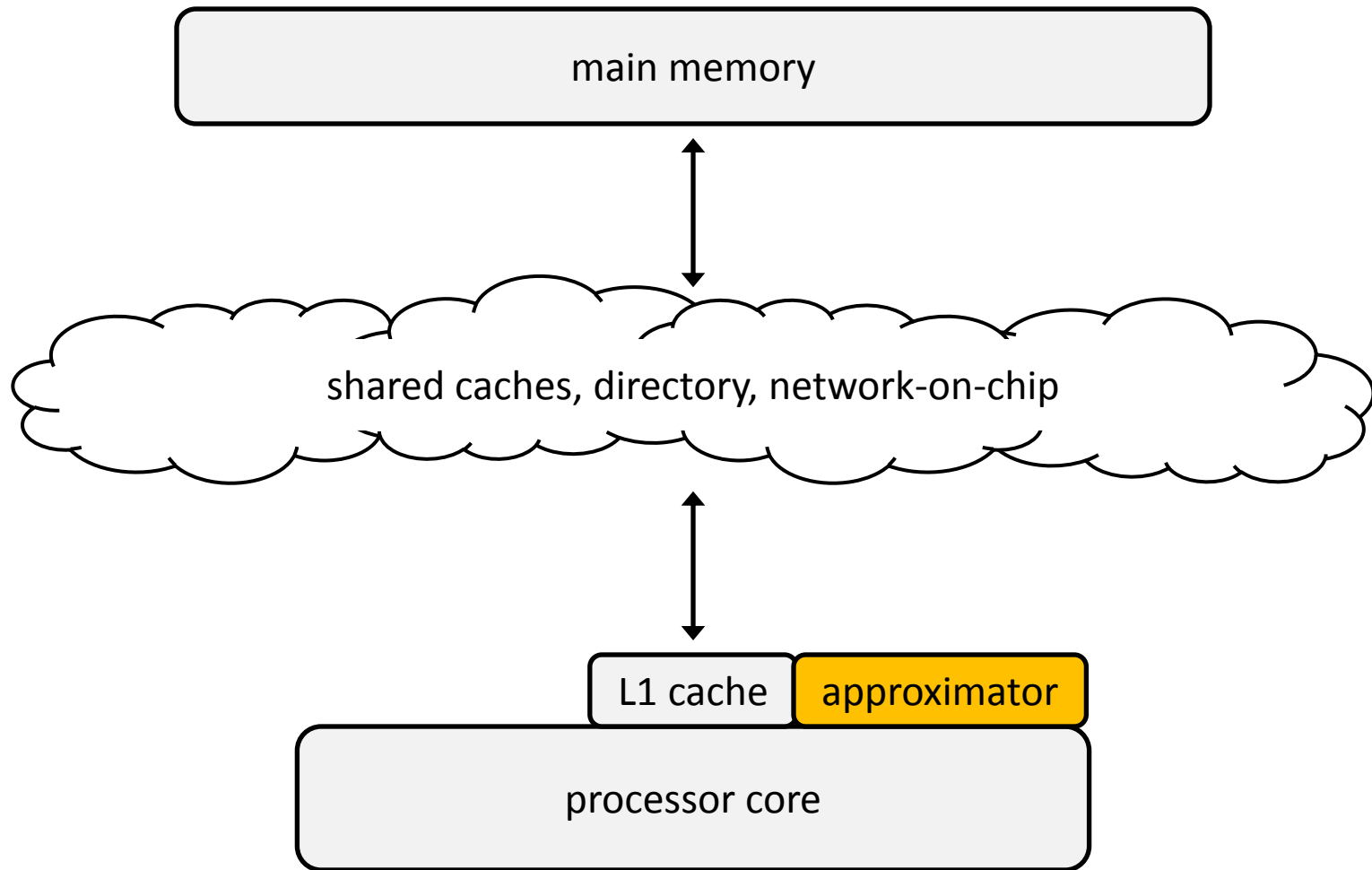
# Outline

- Load Value Approximation
  - Non-Speculative Operation
- Approximator Design
  - Relaxed Confidence Windows
  - Approximation Degree
- Methodology
- Evaluation

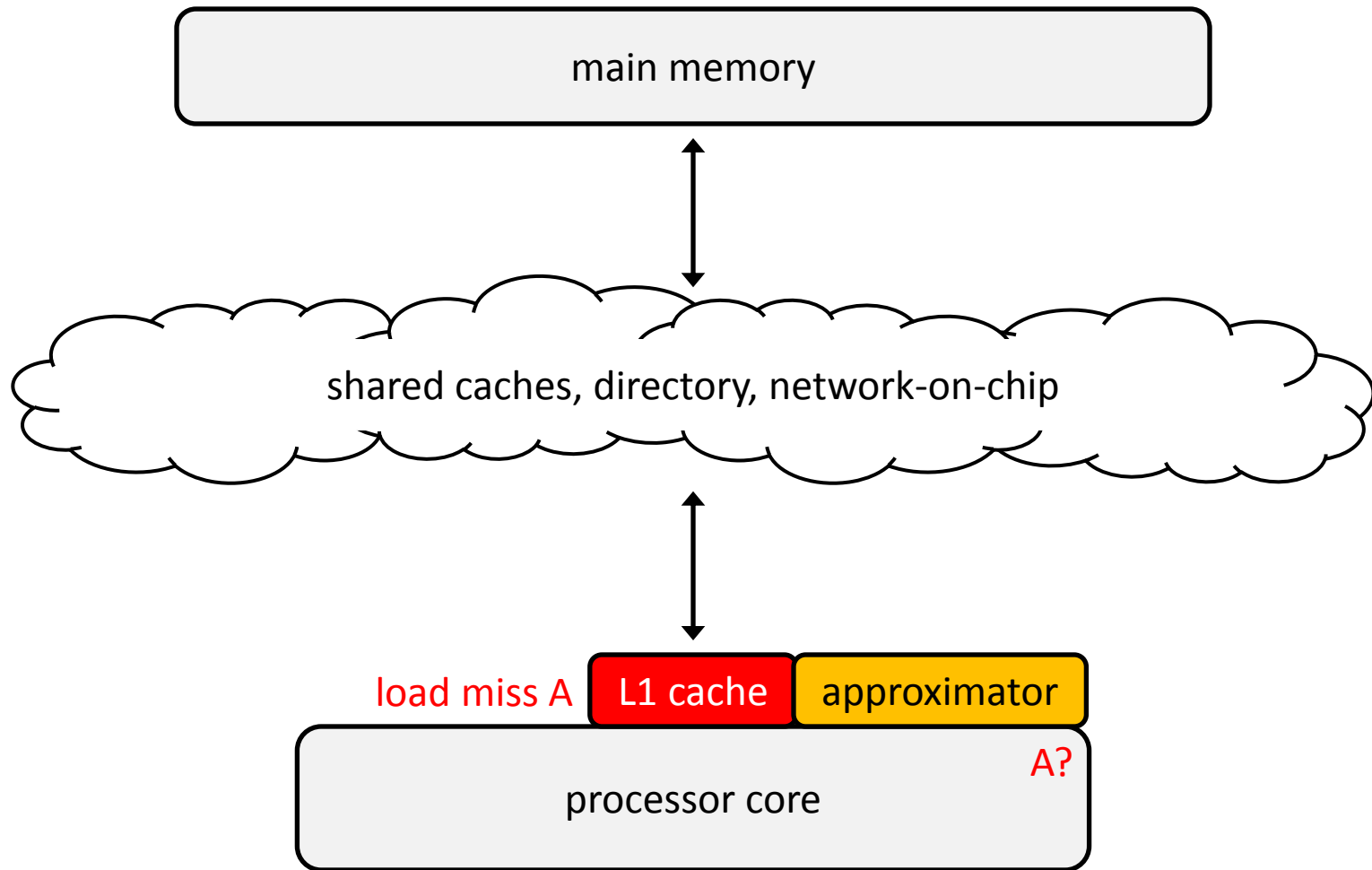
# Load Value Approximation



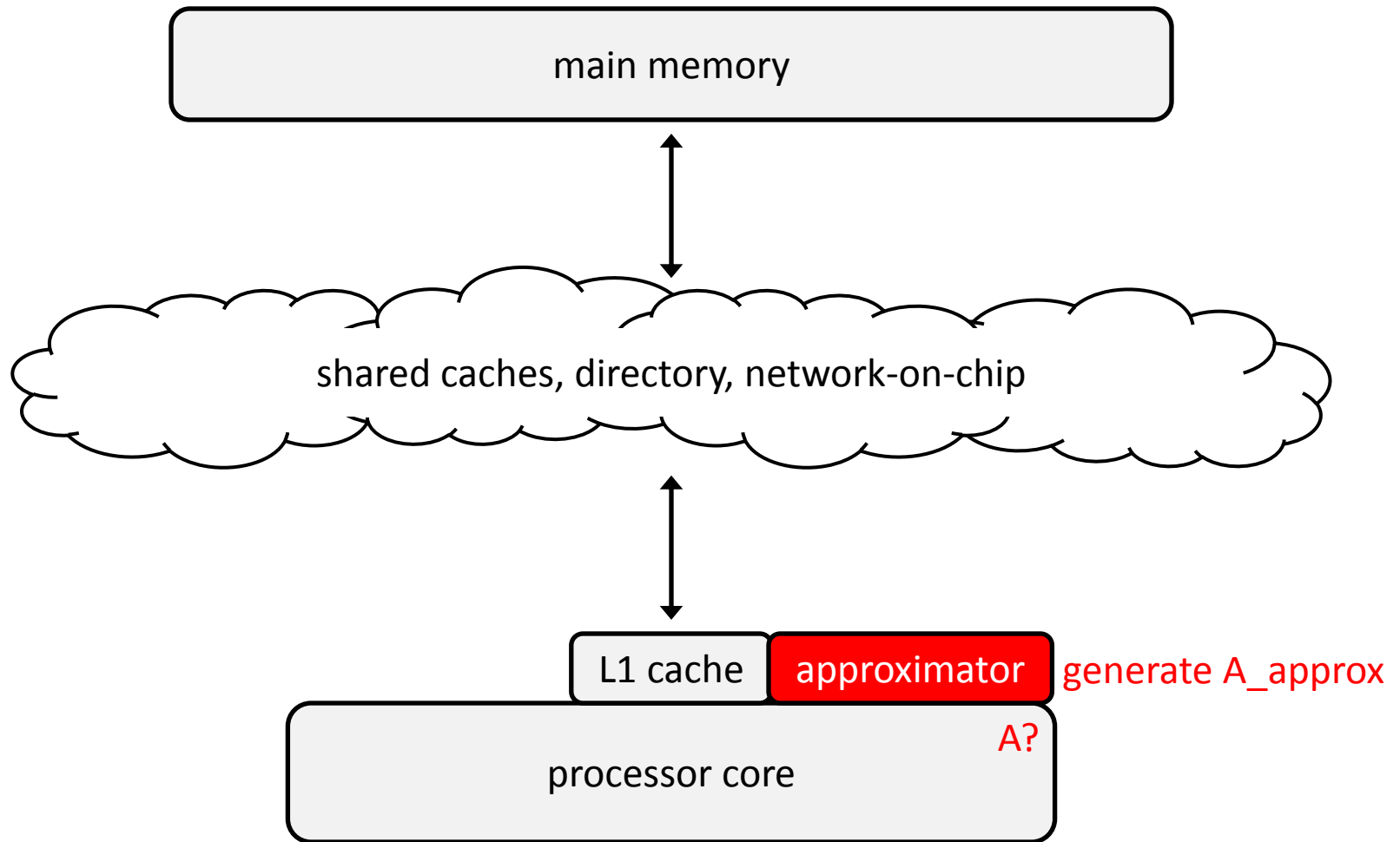
# Load Value Approximation



# Load Value Approximation

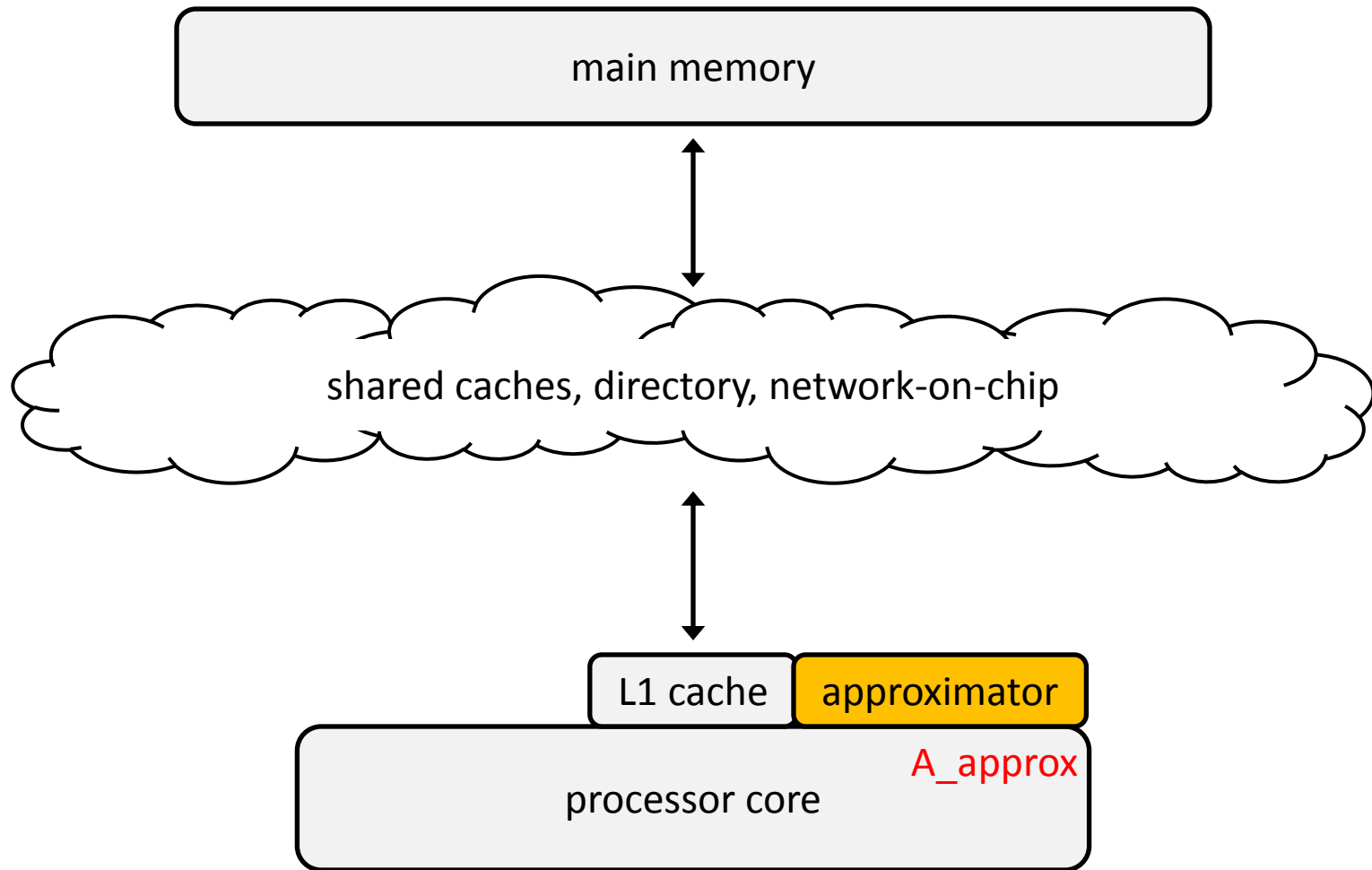


# Load Value Approximation

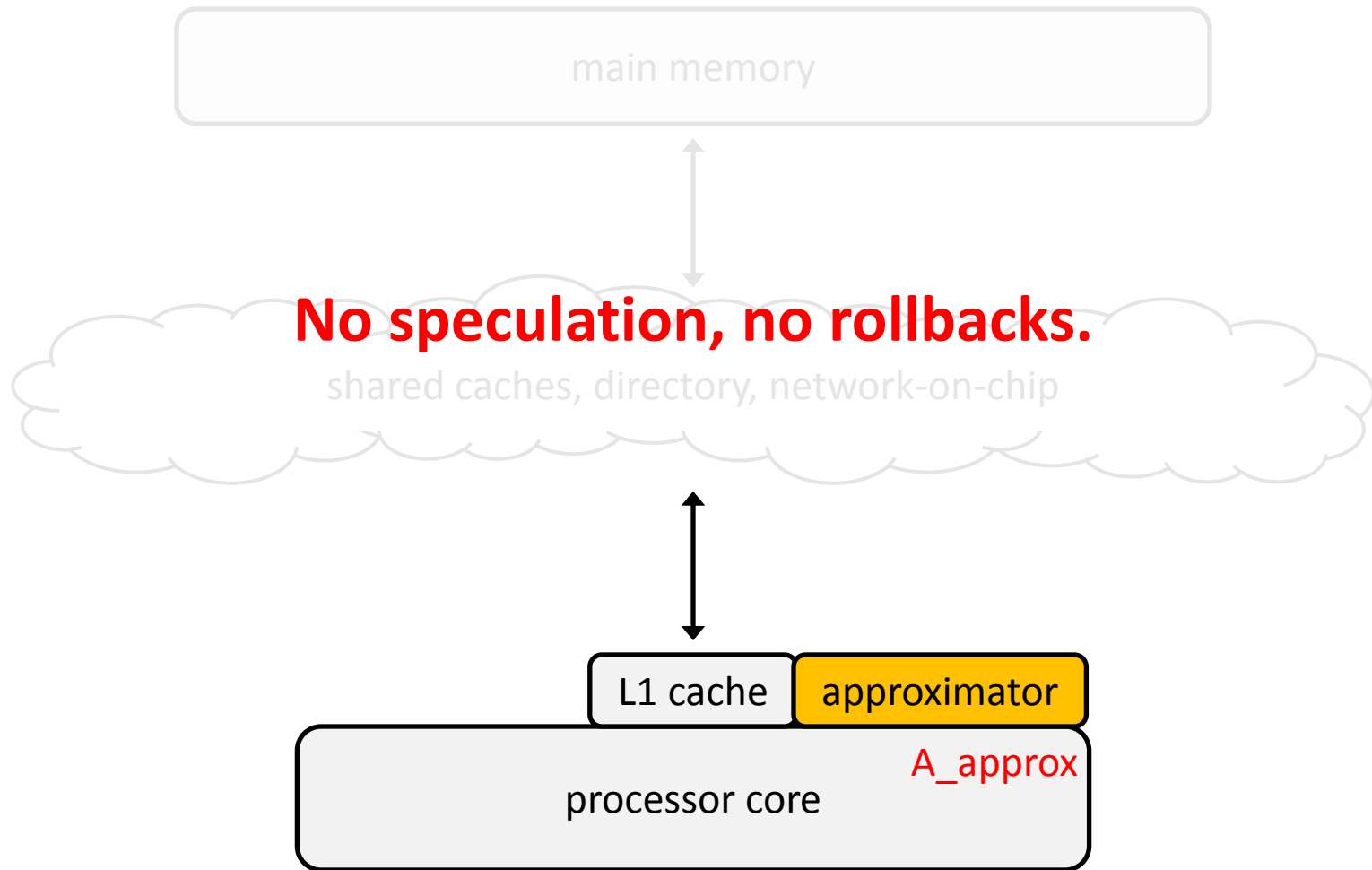




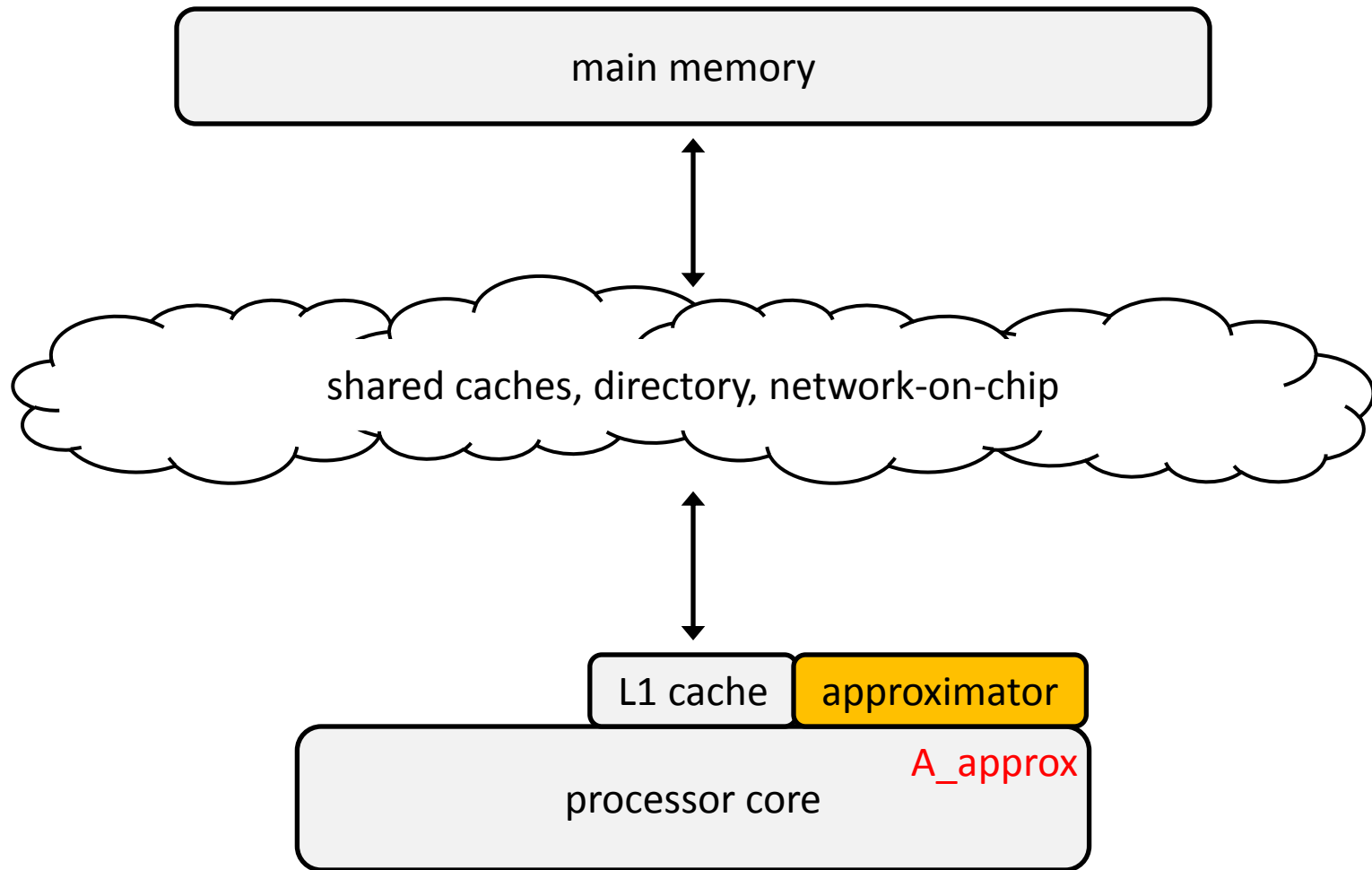
# Load Value Approximation



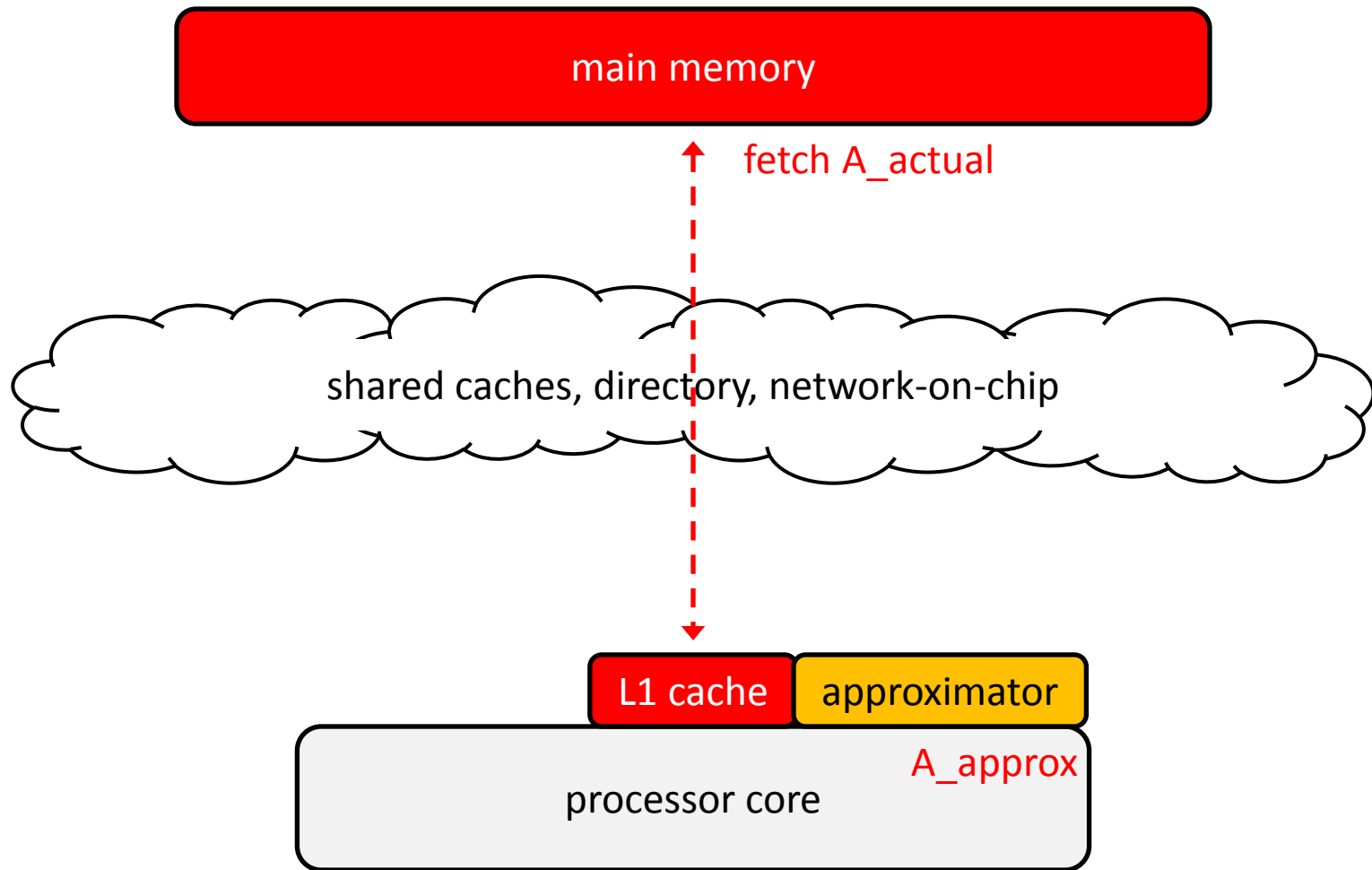
# Load Value Approximation



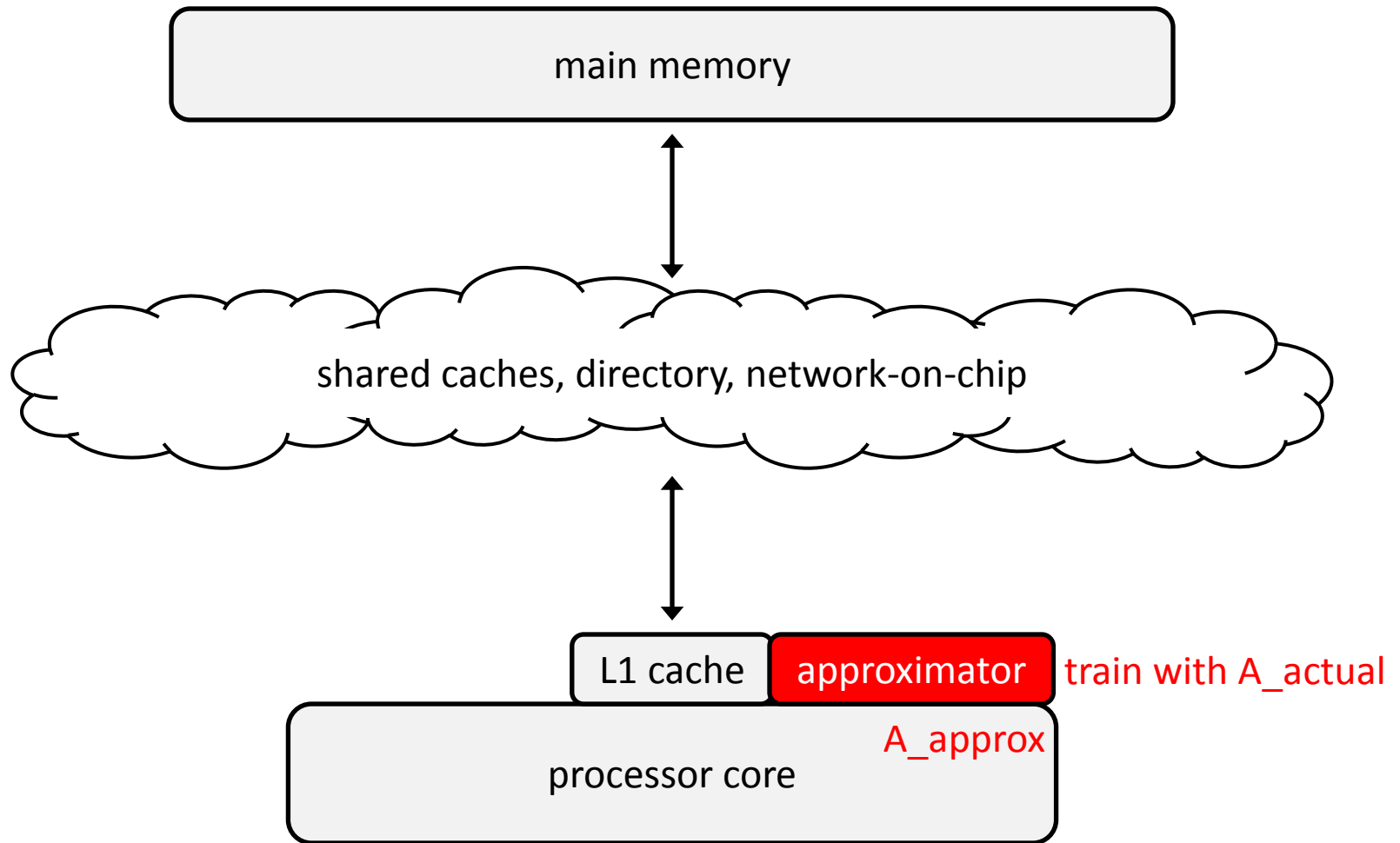
# Load Value Approximation



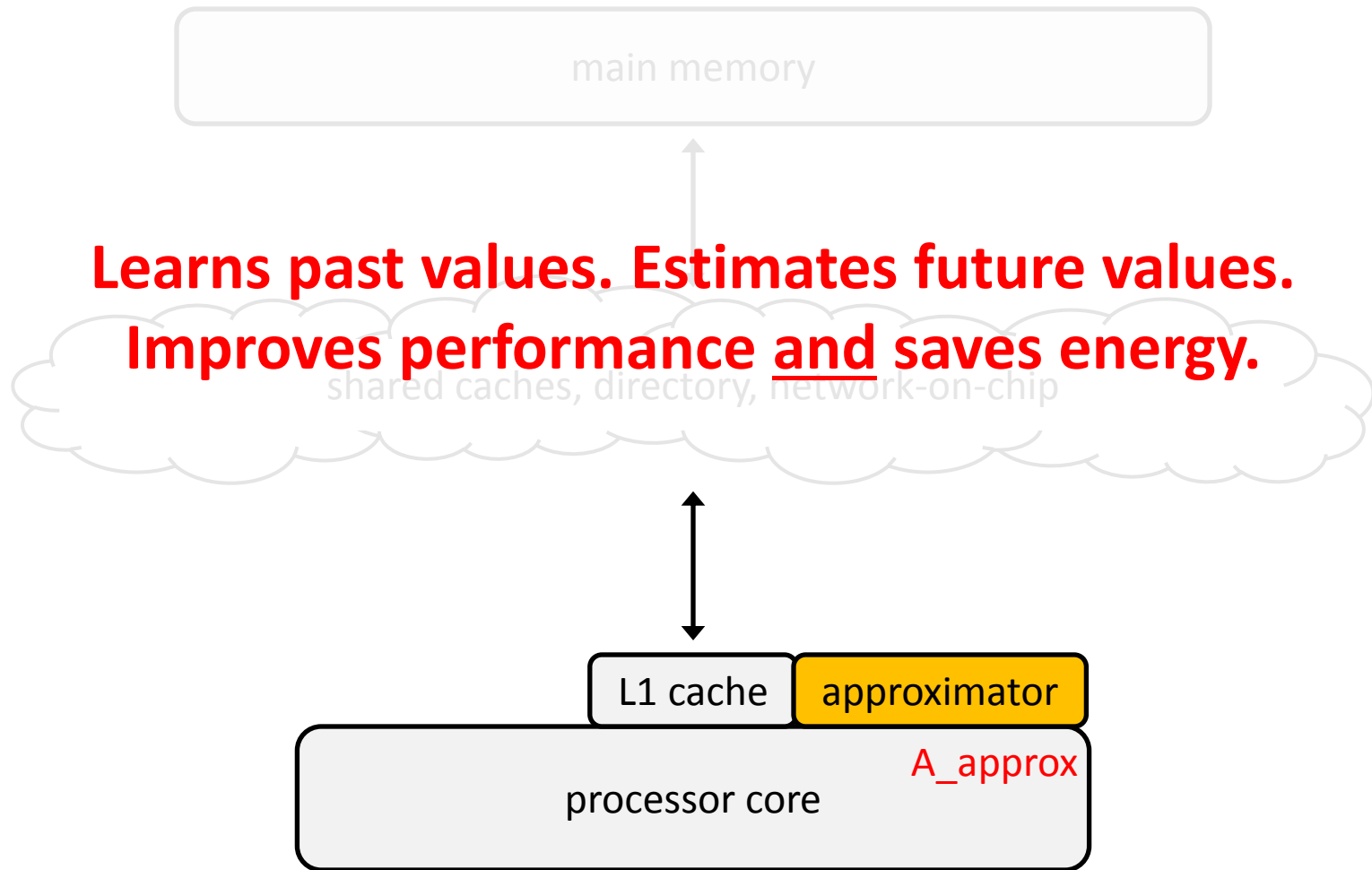
# Load Value Approximation



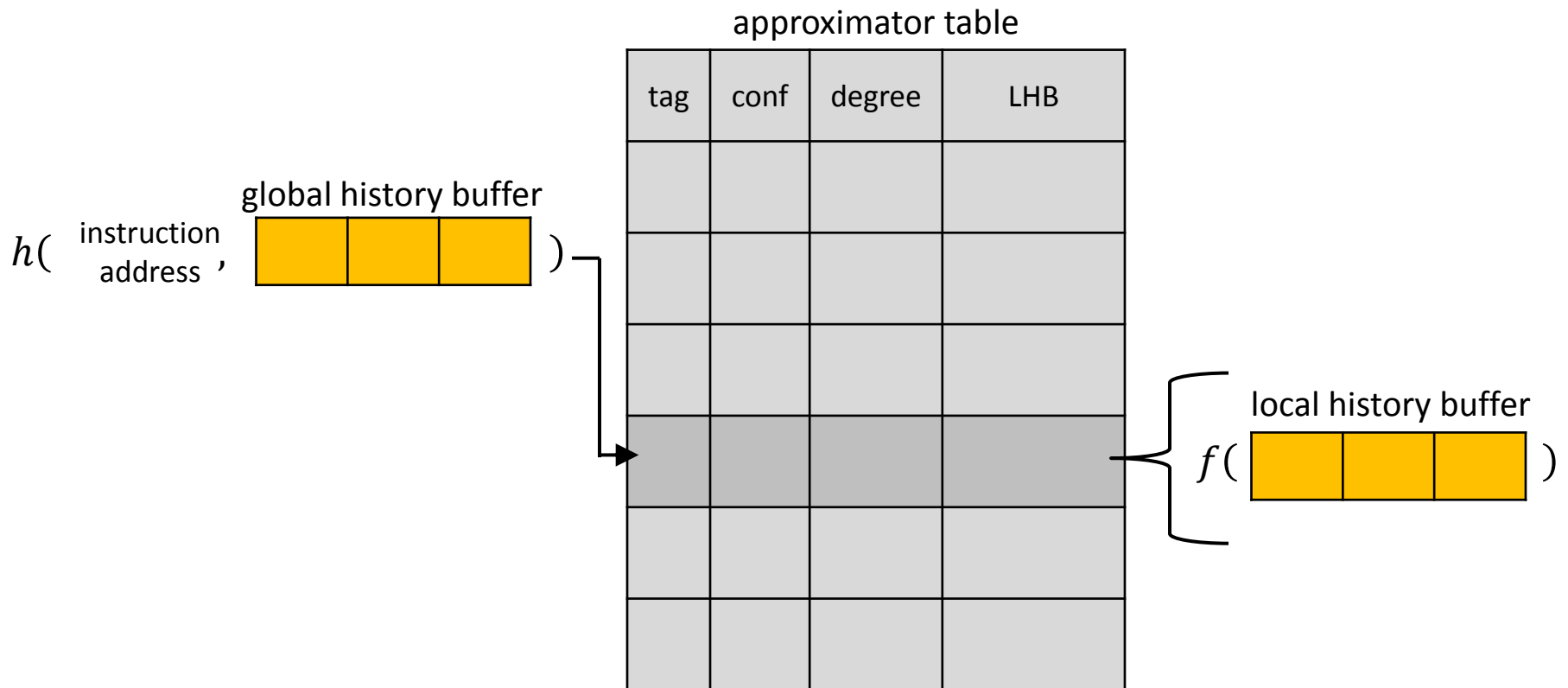
# Load Value Approximation



# Load Value Approximation

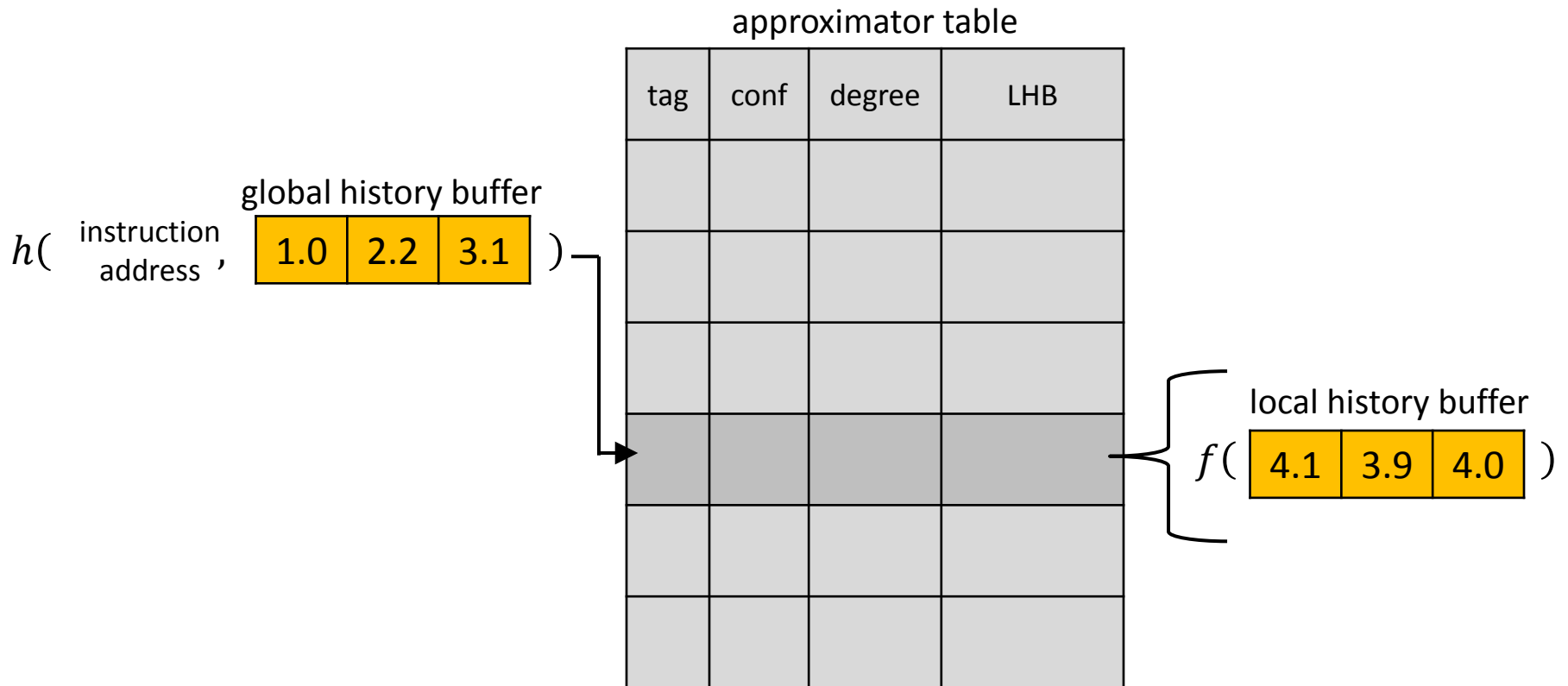


# Approximator Design



# Approximator Design

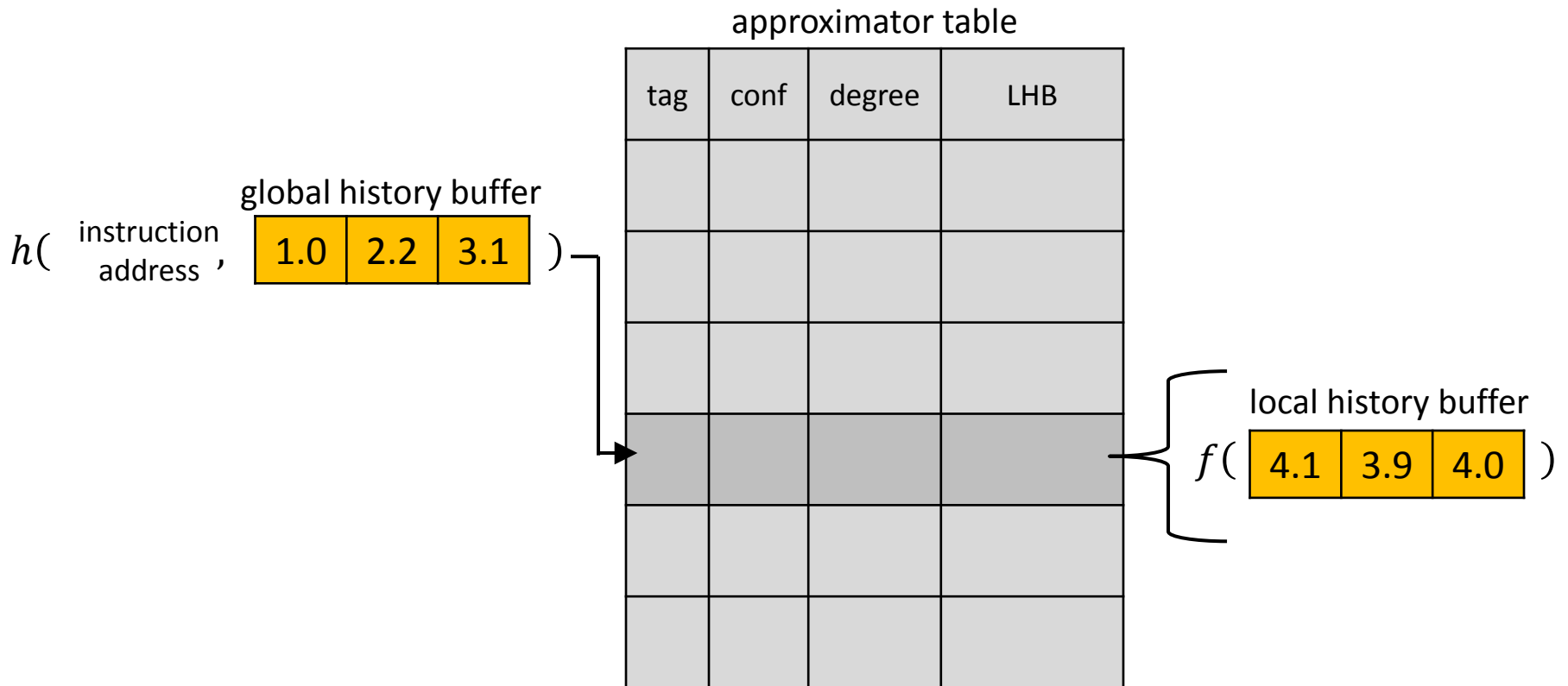
time 





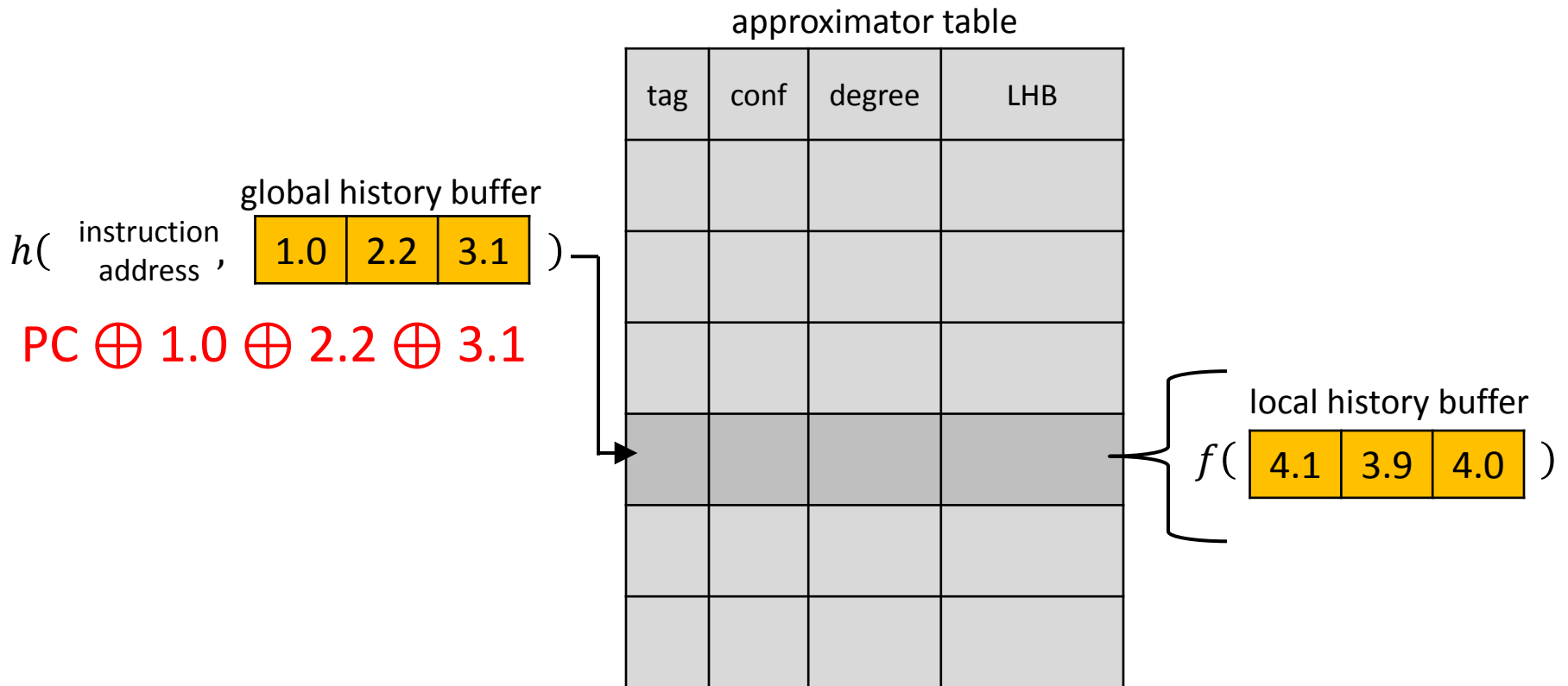
# Approximator Design

load miss A



# Approximator Design

load miss A



# Approximator Design

load miss A



$h(\text{instruction address}, \text{global history buffer})$

1.0 2.2 3.1

$PC \oplus 1.0 \oplus 2.2 \oplus 3.1$

approximator table

tag	conf	degree	LHB

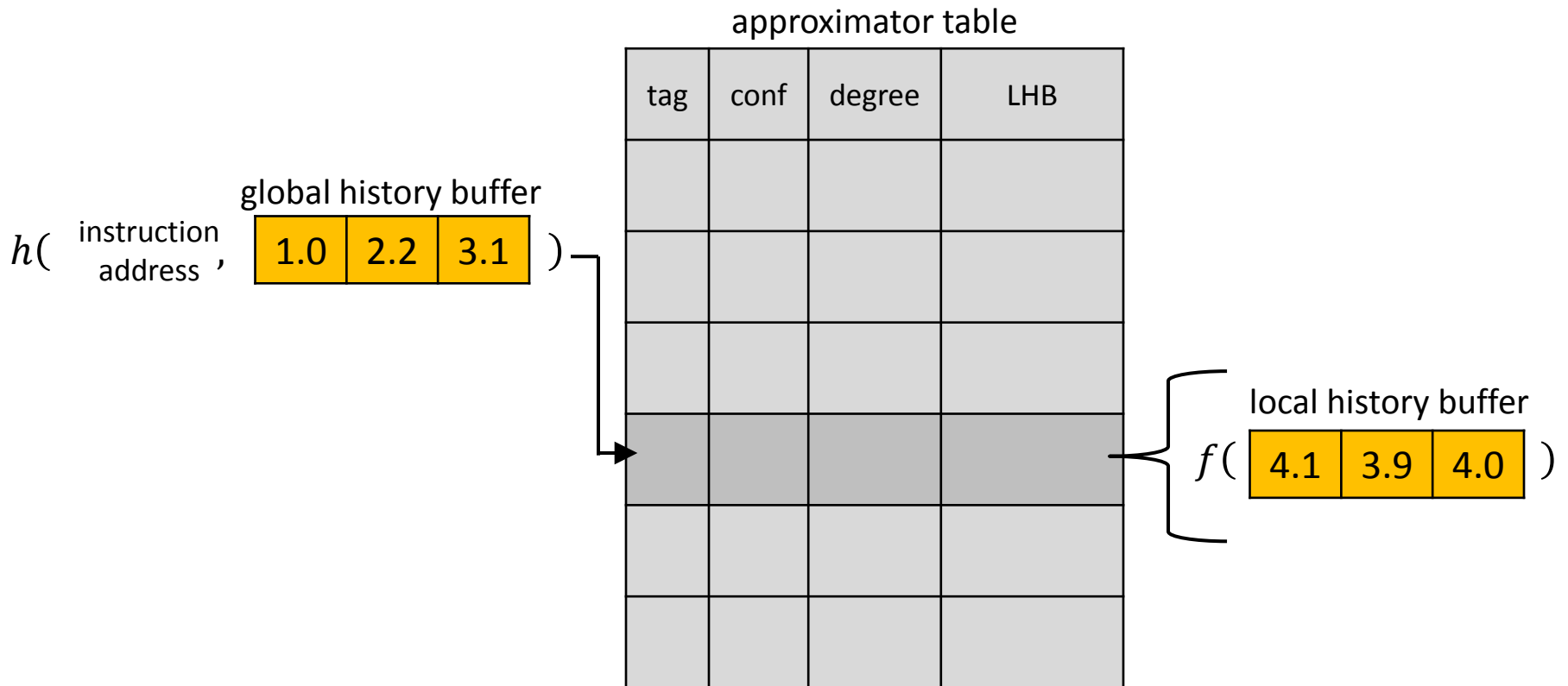
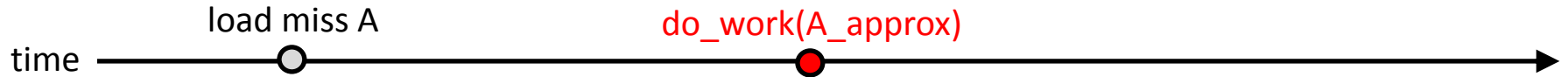
local history buffer

$f(4.1 \ 3.9 \ 4.0)$

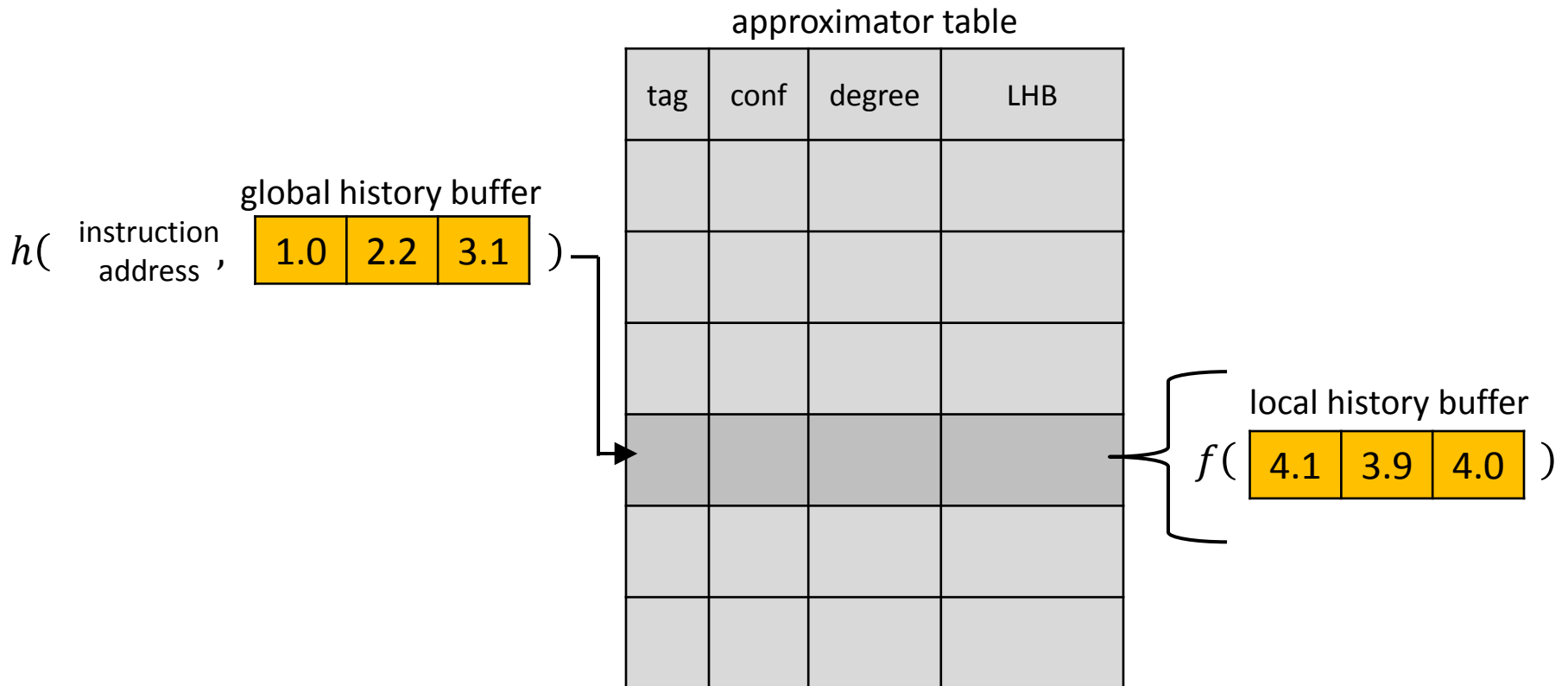
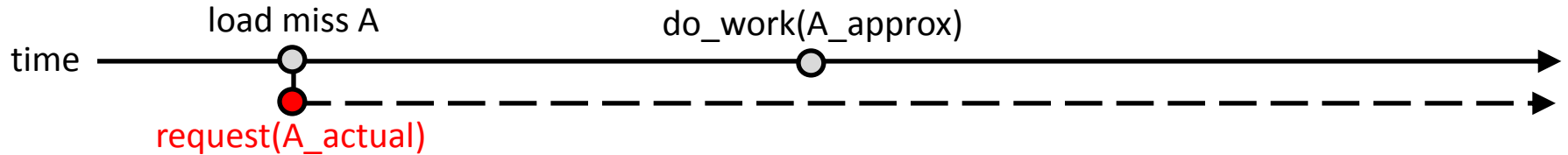
$(4.1 + 3.9 + 4.0) / 3$

$A_{\text{approx}} = 4.0$

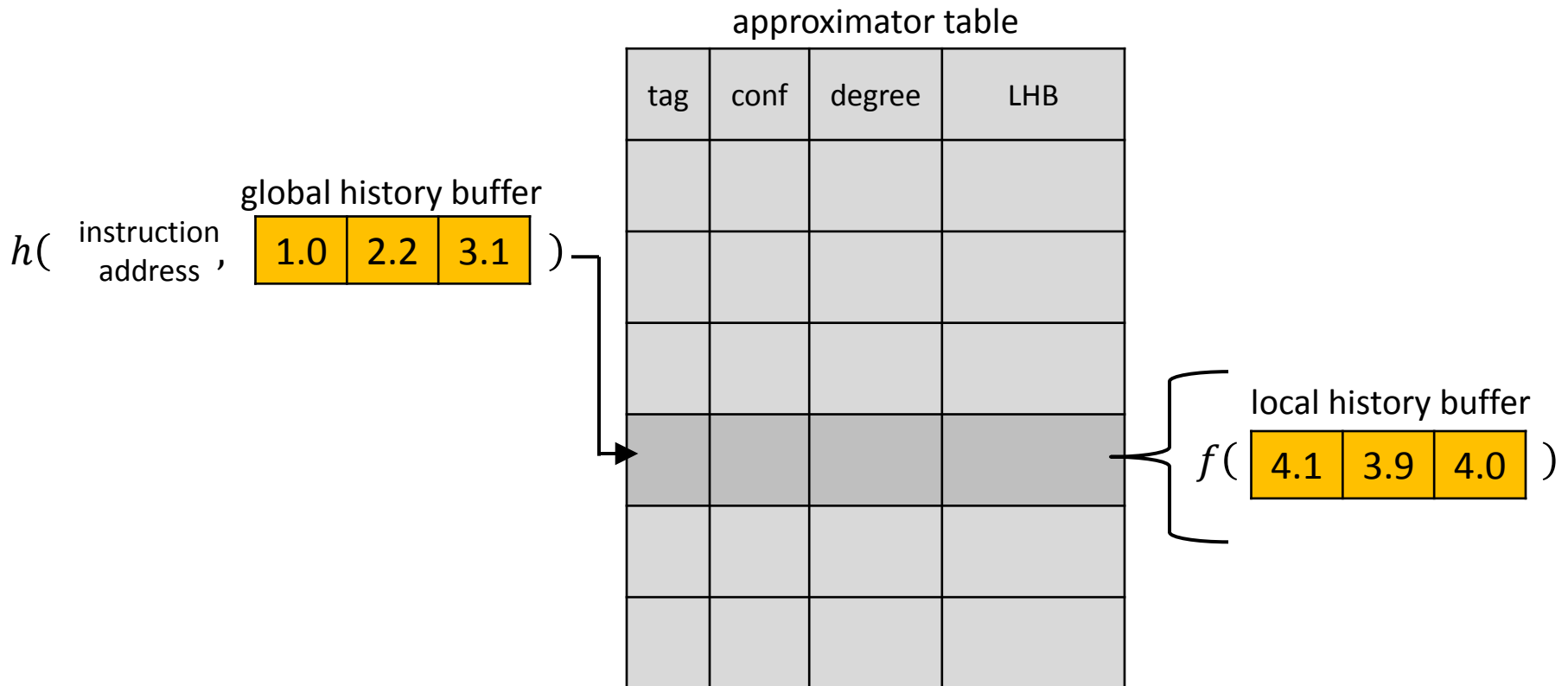
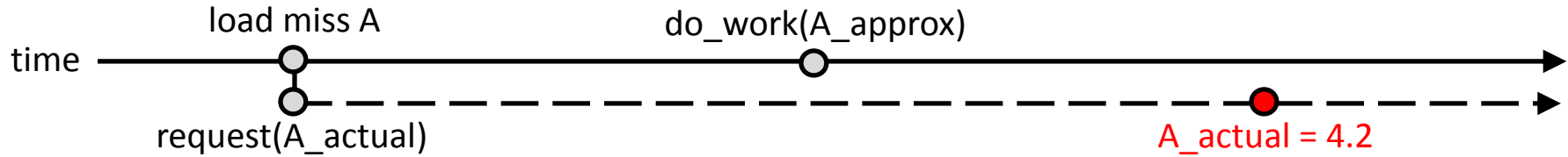
# Approximator Design



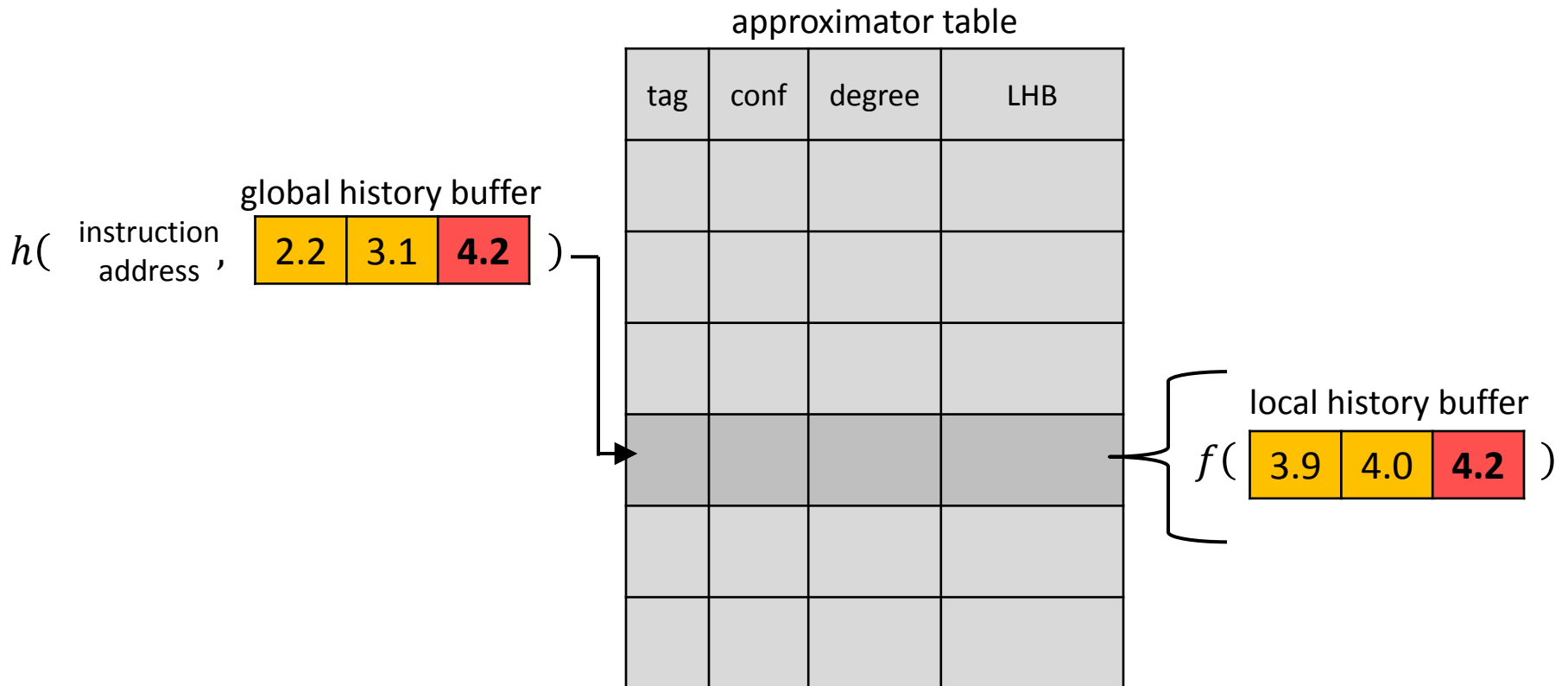
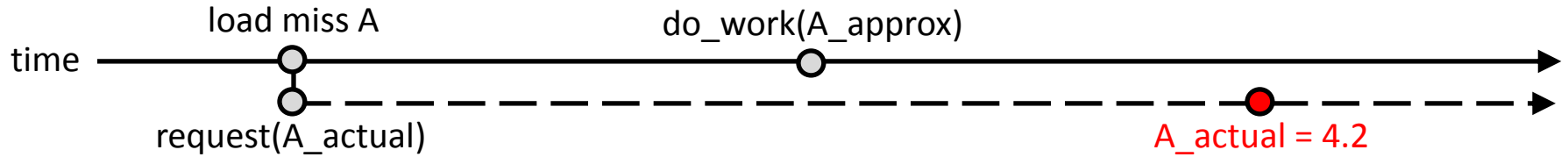
# Approximator Design



# Approximator Design



# Approximator Design



# Approximator Design – Other Considerations

- Floating-point precision
- History buffer sizes
- Stale values

More details in paper.



# Approximator Design

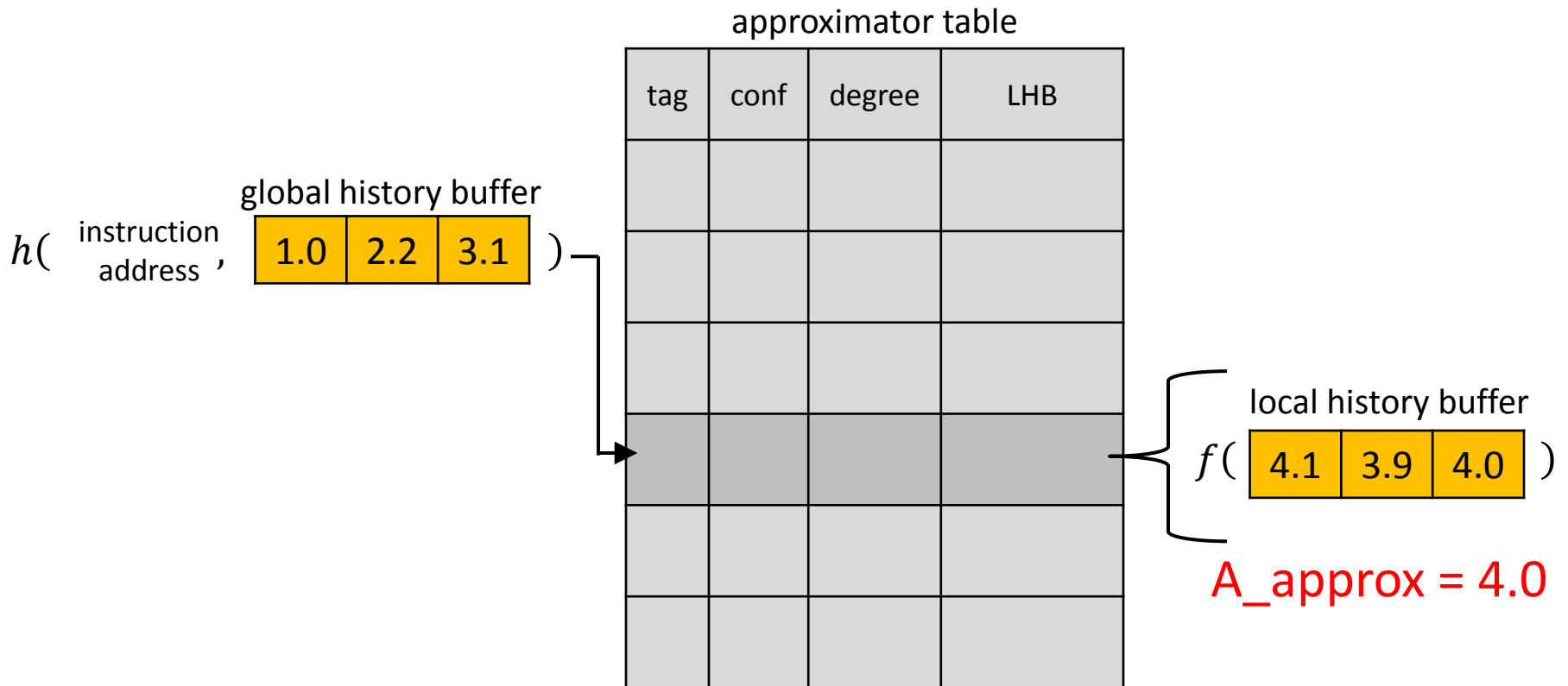
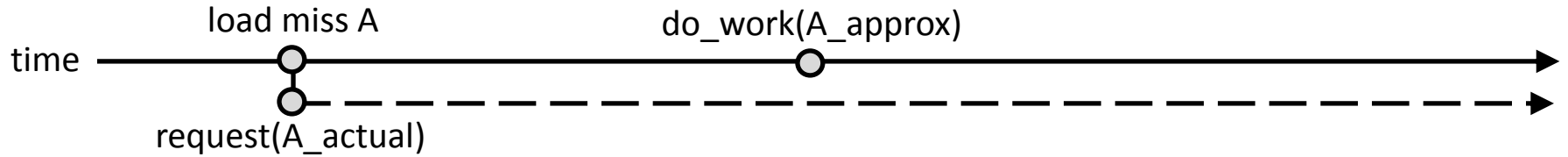
## Relaxed Confidence Windows

- How do we avoid making bad approximations?
- Trade-off performance and error.

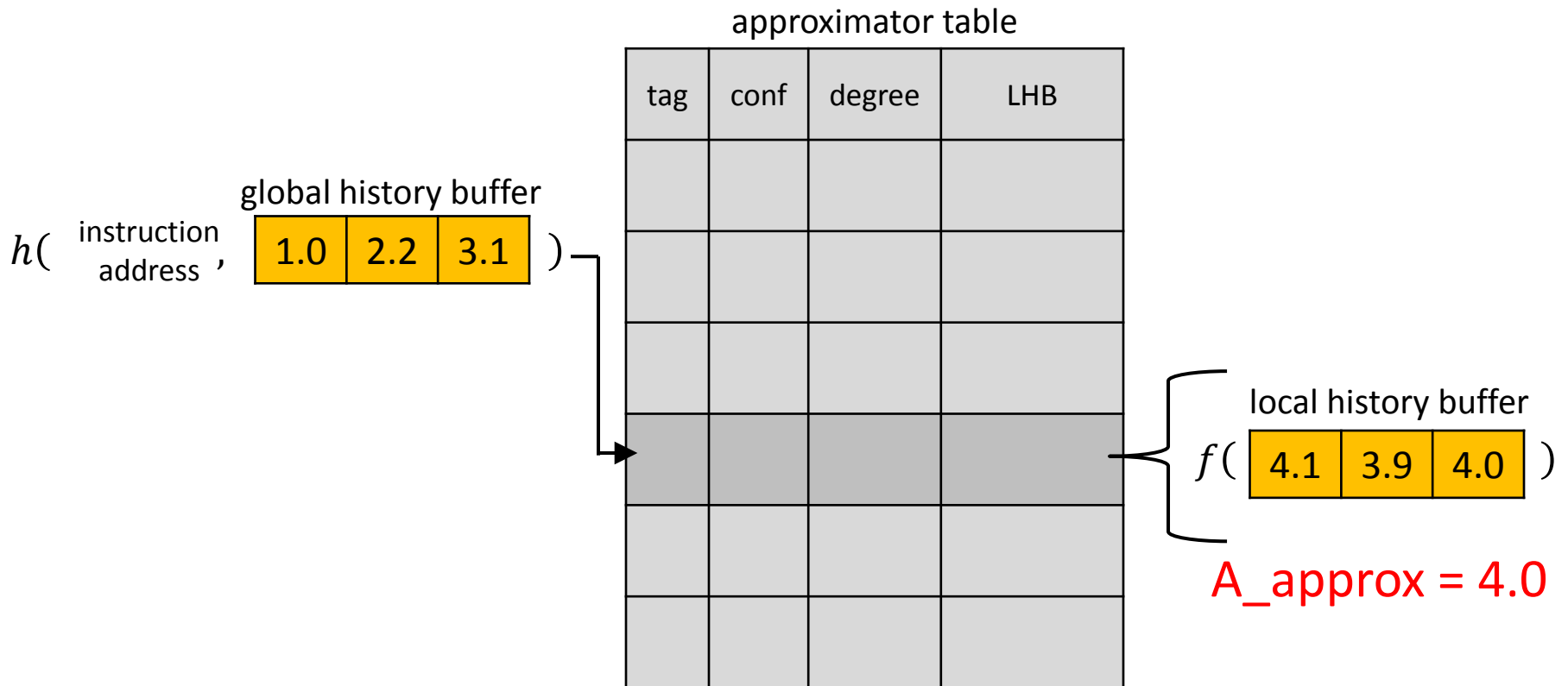
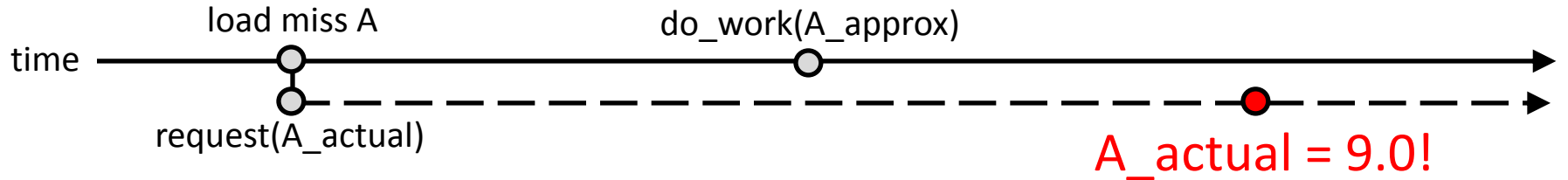
## Approximation Degree

- Do we need to fetch the actual value from memory every time?
- Trade-off energy and error.

# Relaxed Confidence Windows



# Relaxed Confidence Windows



# Relaxed Confidence Windows

tag	conf	degree	LHB
-----	------	--------	-----

When approximating:

if  $conf \geq 0$ : use  $A_{approx}$

else: don't use  $A_{approx}$

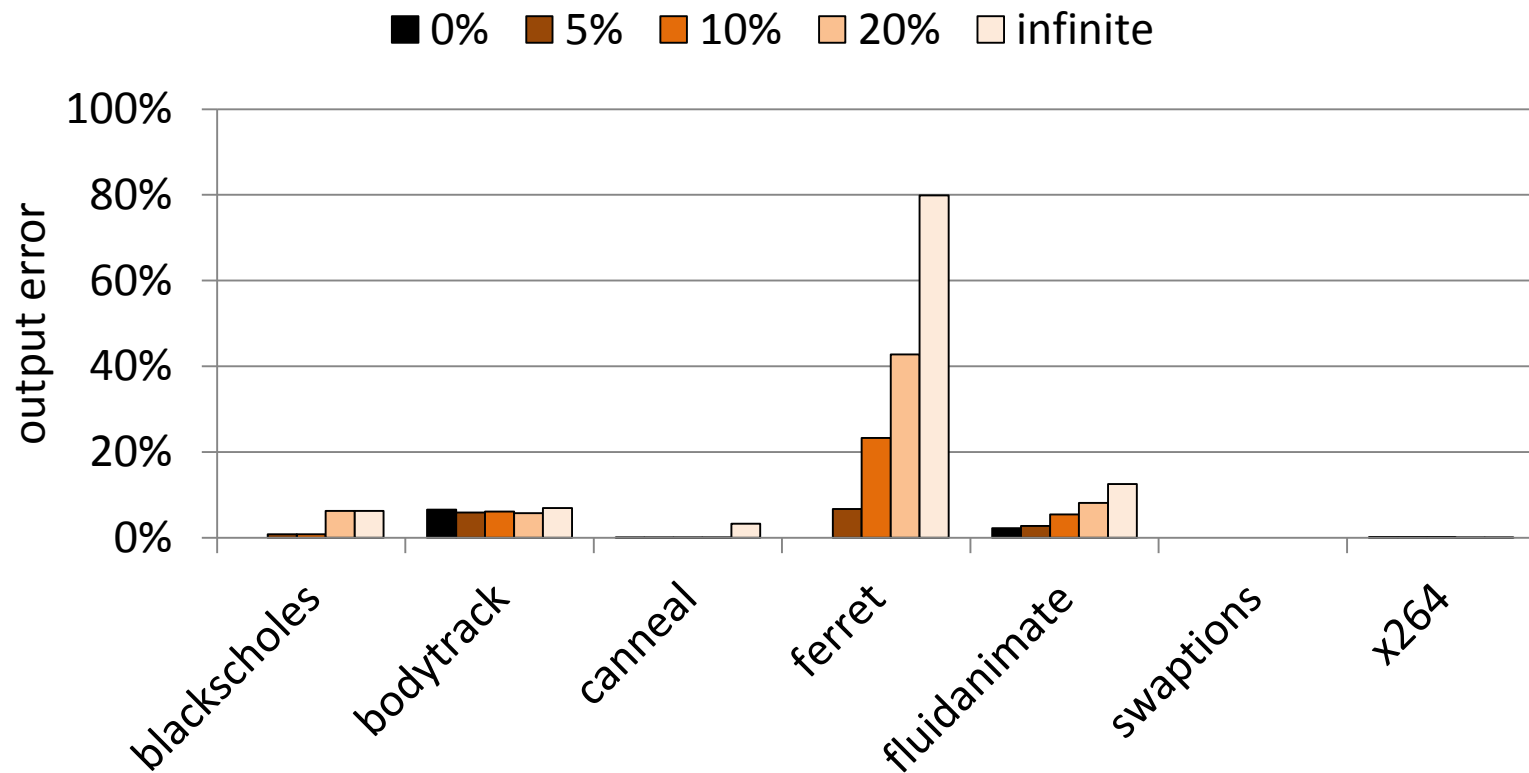
When updating:

if  $A_{approx}$ ,  $A_{actual}$  differ by  $\leq \mathbf{CONF\_WINDOW\%}$ :  $conf++$

else:  $conf--$

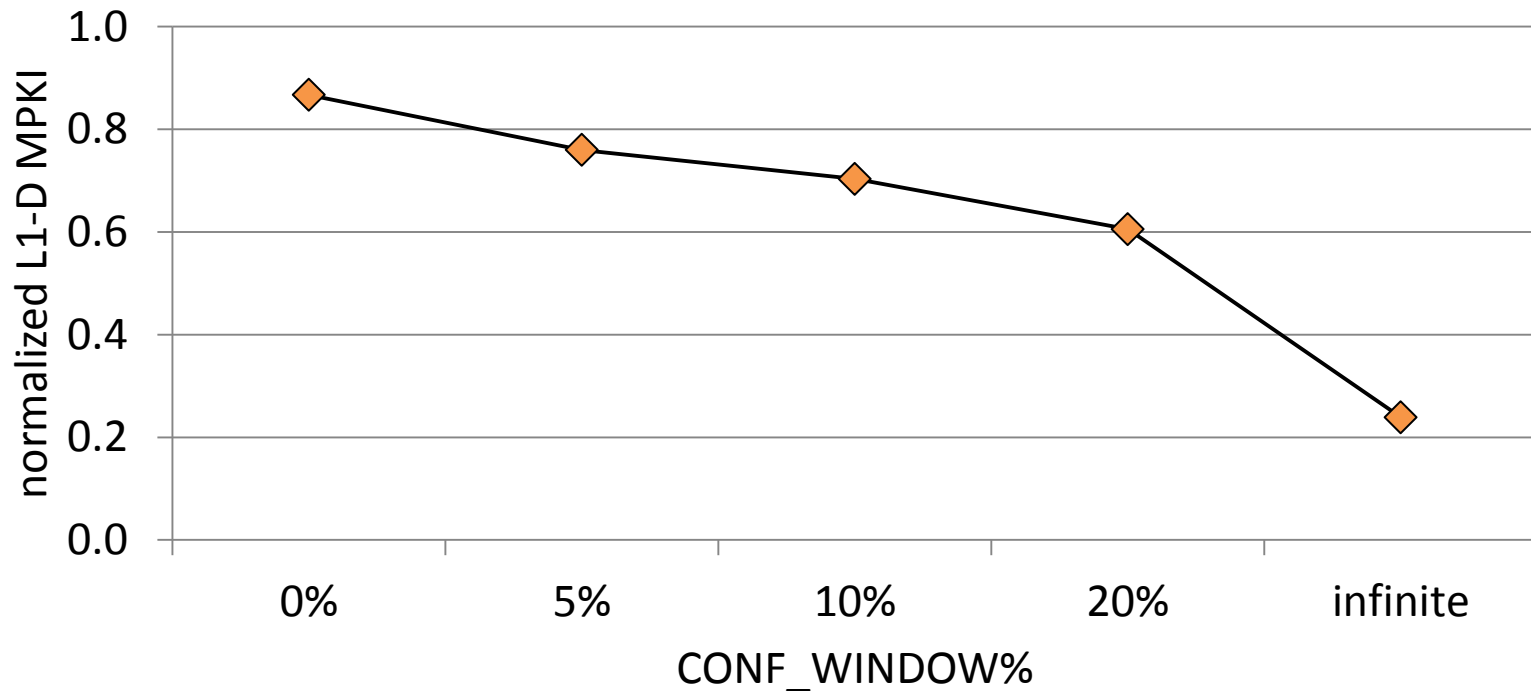
# Relaxed Confidence Windows – Output Error

Varying *CONF\_WINDOW*%:



# Relaxed Confidence Windows – L1-D MPKI

Varying *CONF\_WINDOW%*:



# Approximator Design

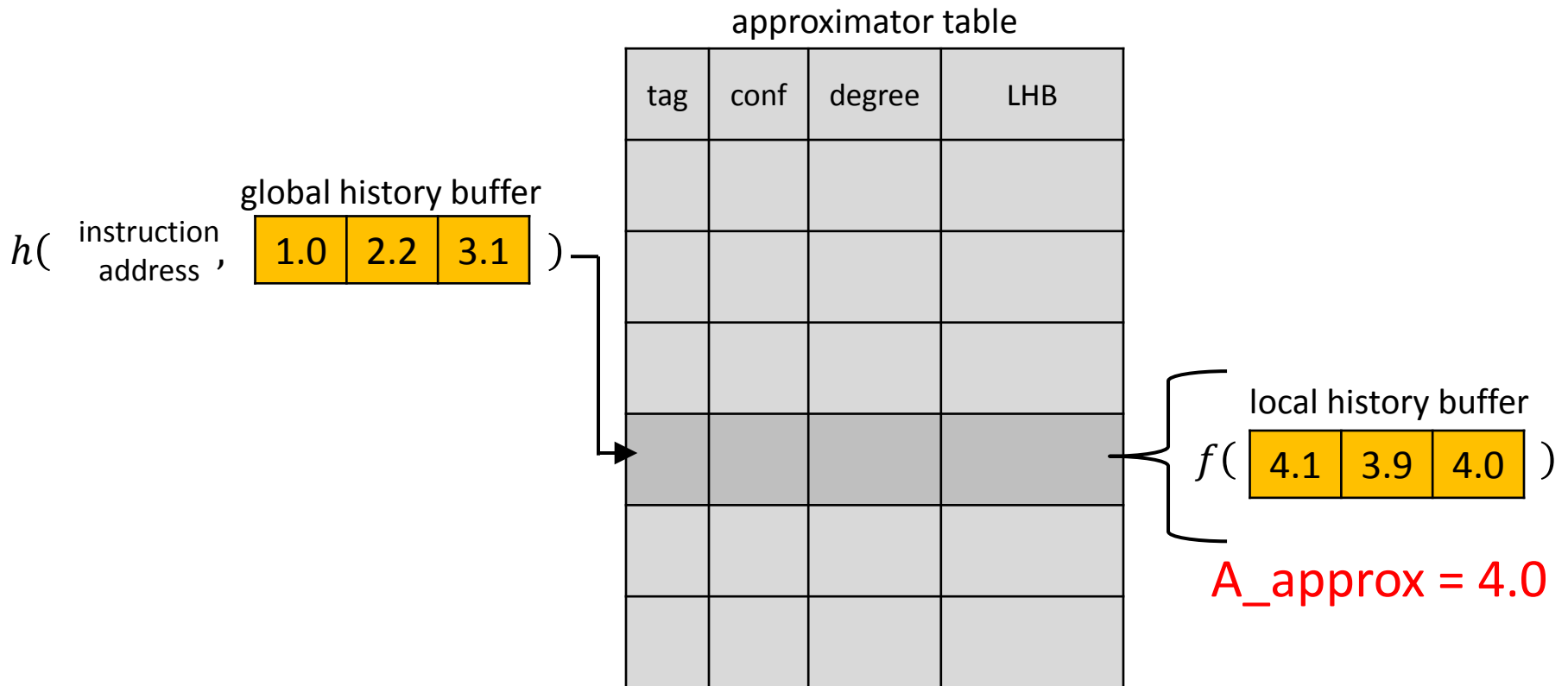
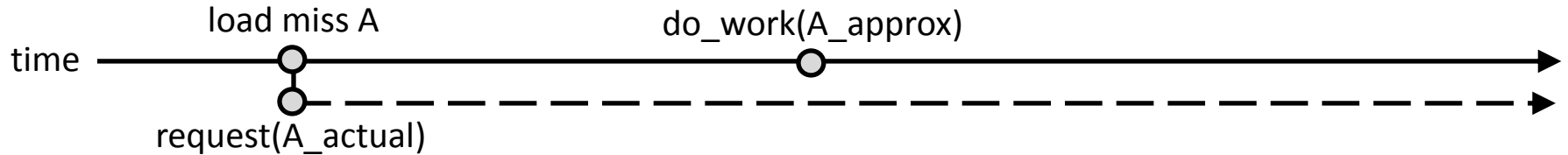
## Relaxed Confidence Windows

- How do we avoid making bad approximations?
- Trade-off performance and error.

## Approximation Degree

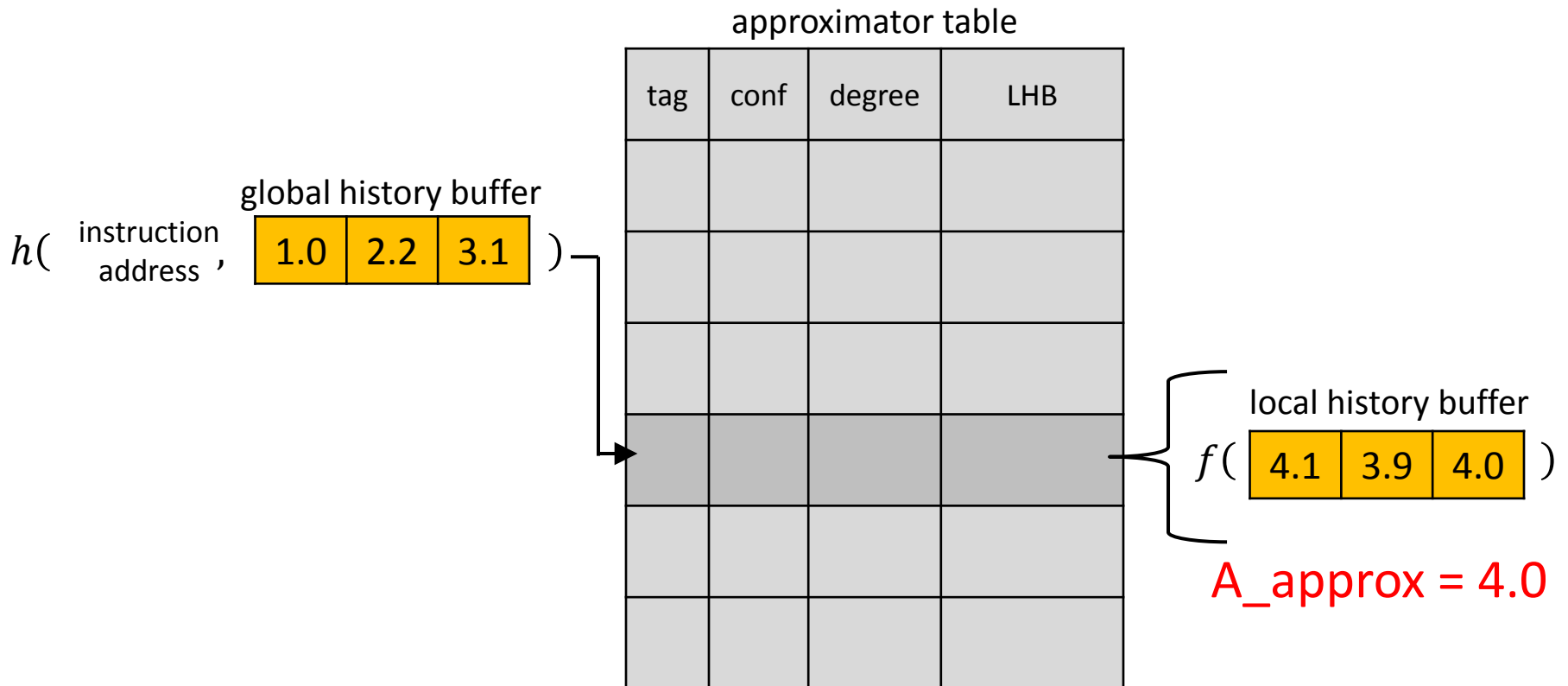
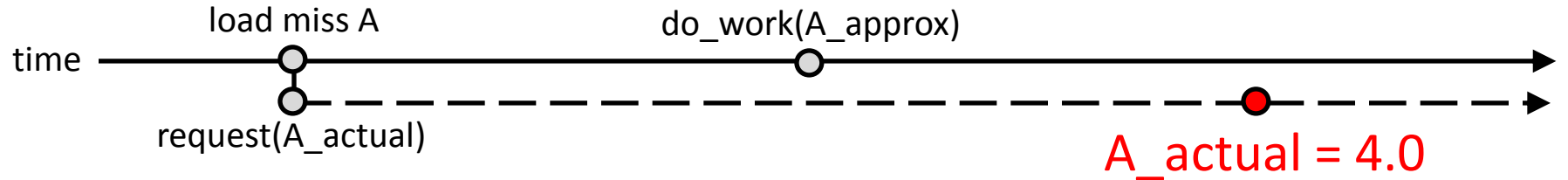
- **Do we need to fetch the actual value from memory every time?**
- **Trade-off energy and error.**

# Approximation Degree

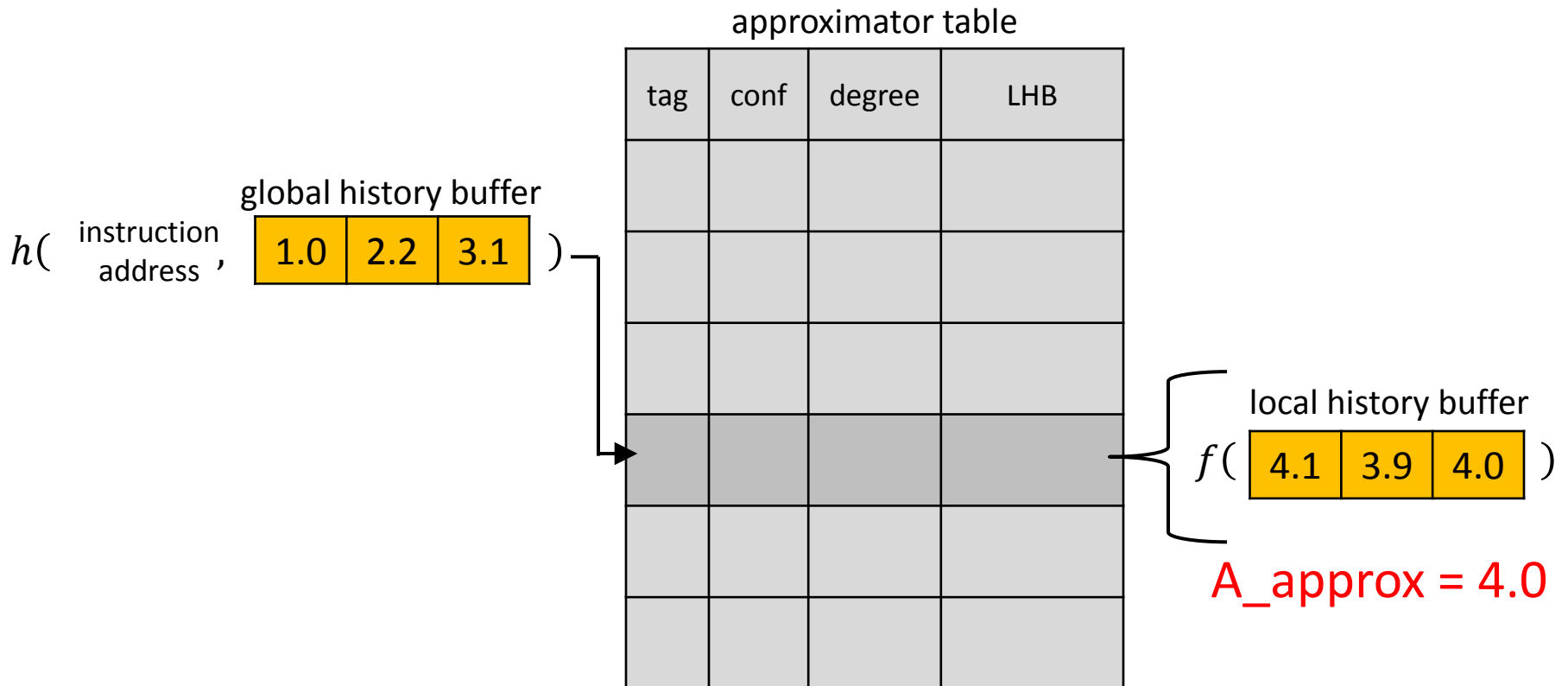
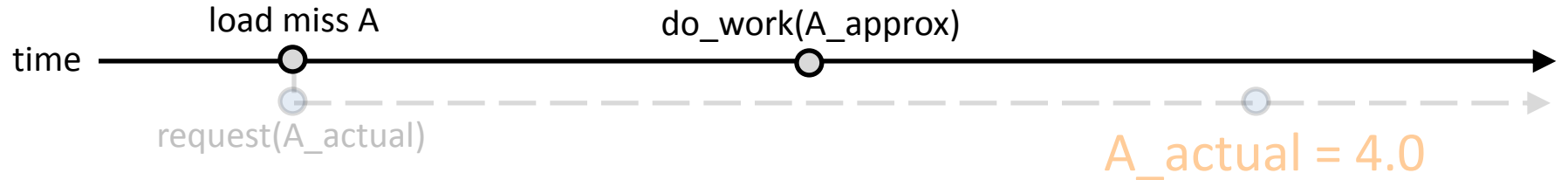




# Approximation Degree



# Approximation Degree



# Approximation Degree

tag	conf	degree	LHB
-----	------	--------	-----

When approximating:

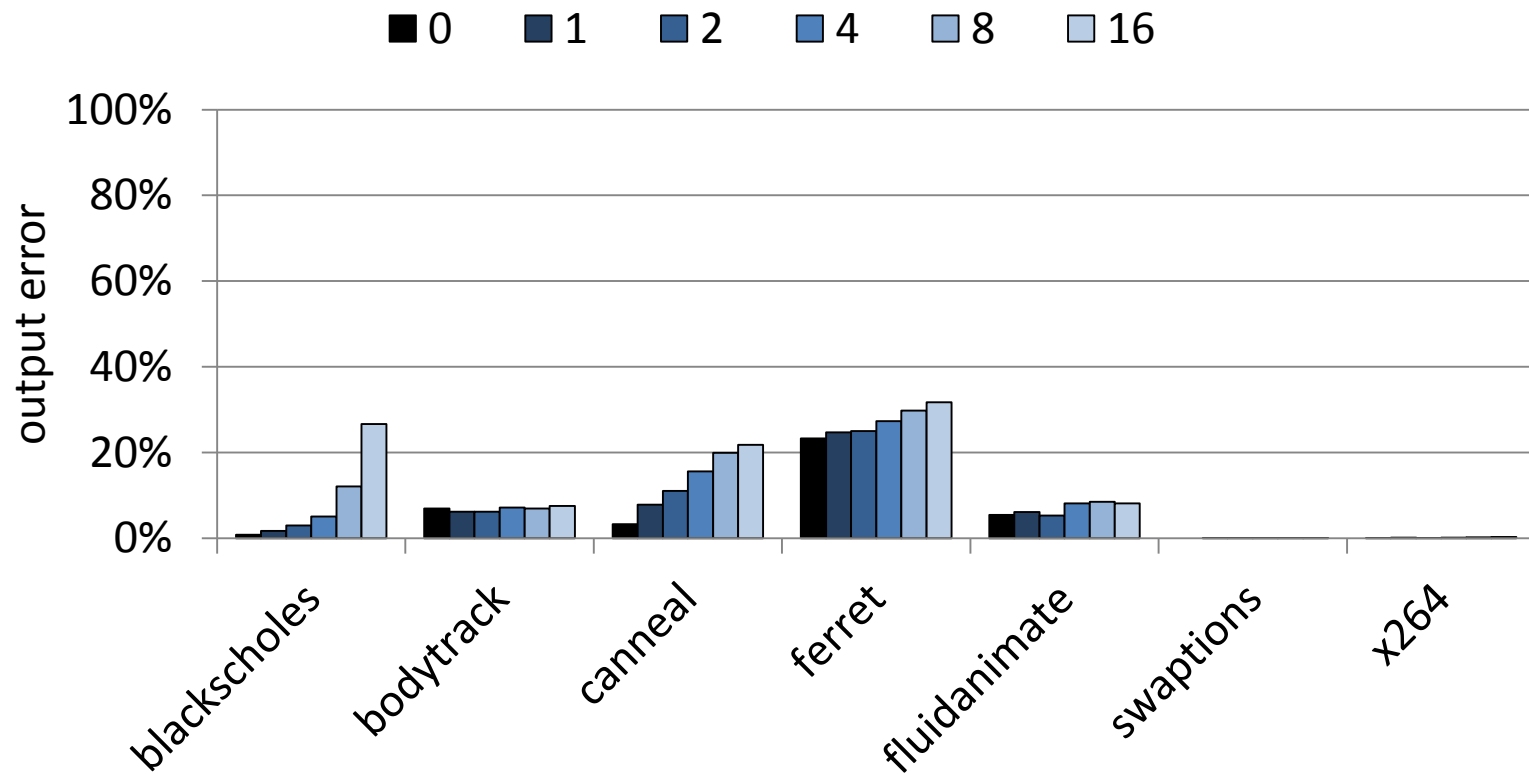
if *degree* == **APPROX\_DEGREE**: fetch *A\_actual*  
else: don't fetch *A\_actual*

When updating:

if *degree* == **APPROX\_DEGREE**: *degree* = 0  
else: *degree*++

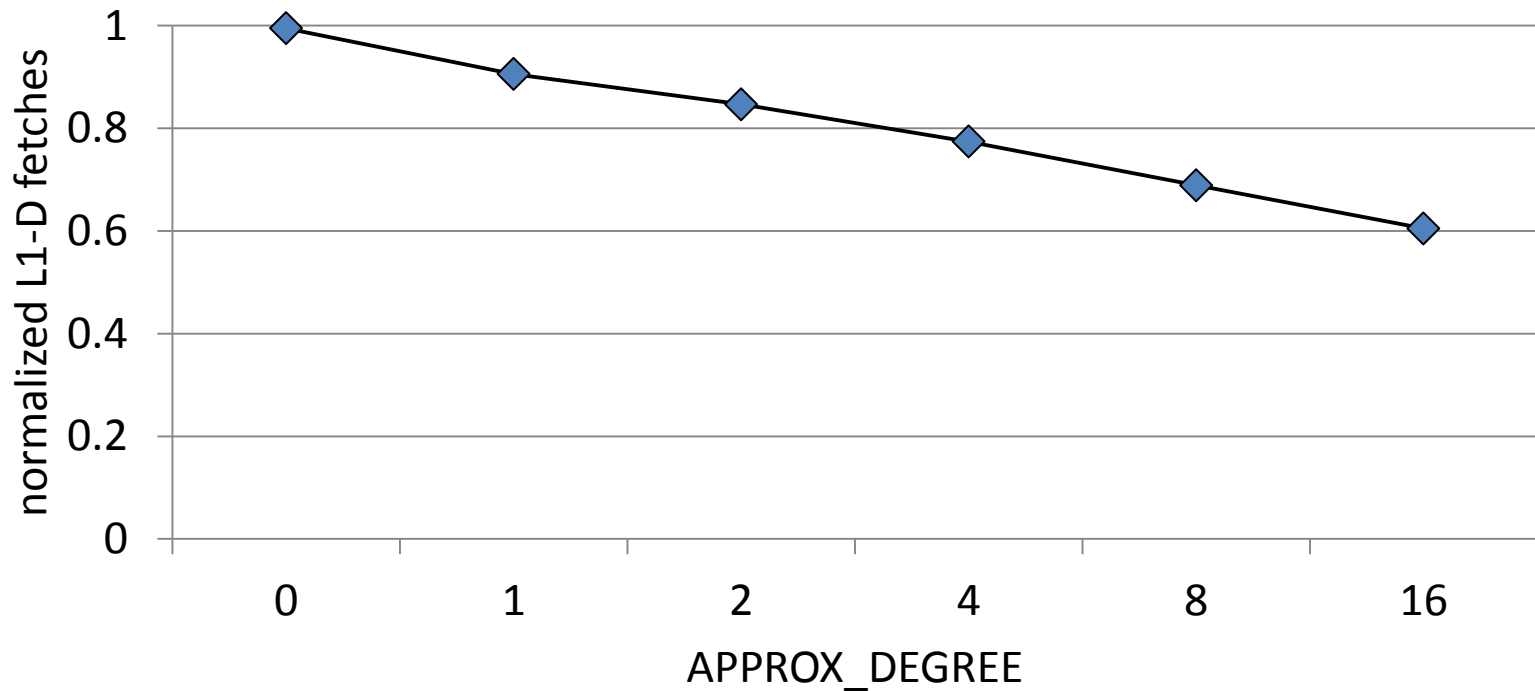
# Approximation Degree – Output Error

Varying *APPROX\_DEGREE*:



# Approximation Degree – L1-D Fetches

Varying *APPROX\_DEGREE*:



# Methodology

## Multi-threaded approximate applications

- PARSEC benchmark suite [Bienia, Princeton 2011]
- Programmer annotations and ISA extensions [Esmaeilzadeh, ASPLOS 2012]

## Approximator design space exploration

- Pin dynamic binary instrumentation tool [Luk, PLDI 2005]

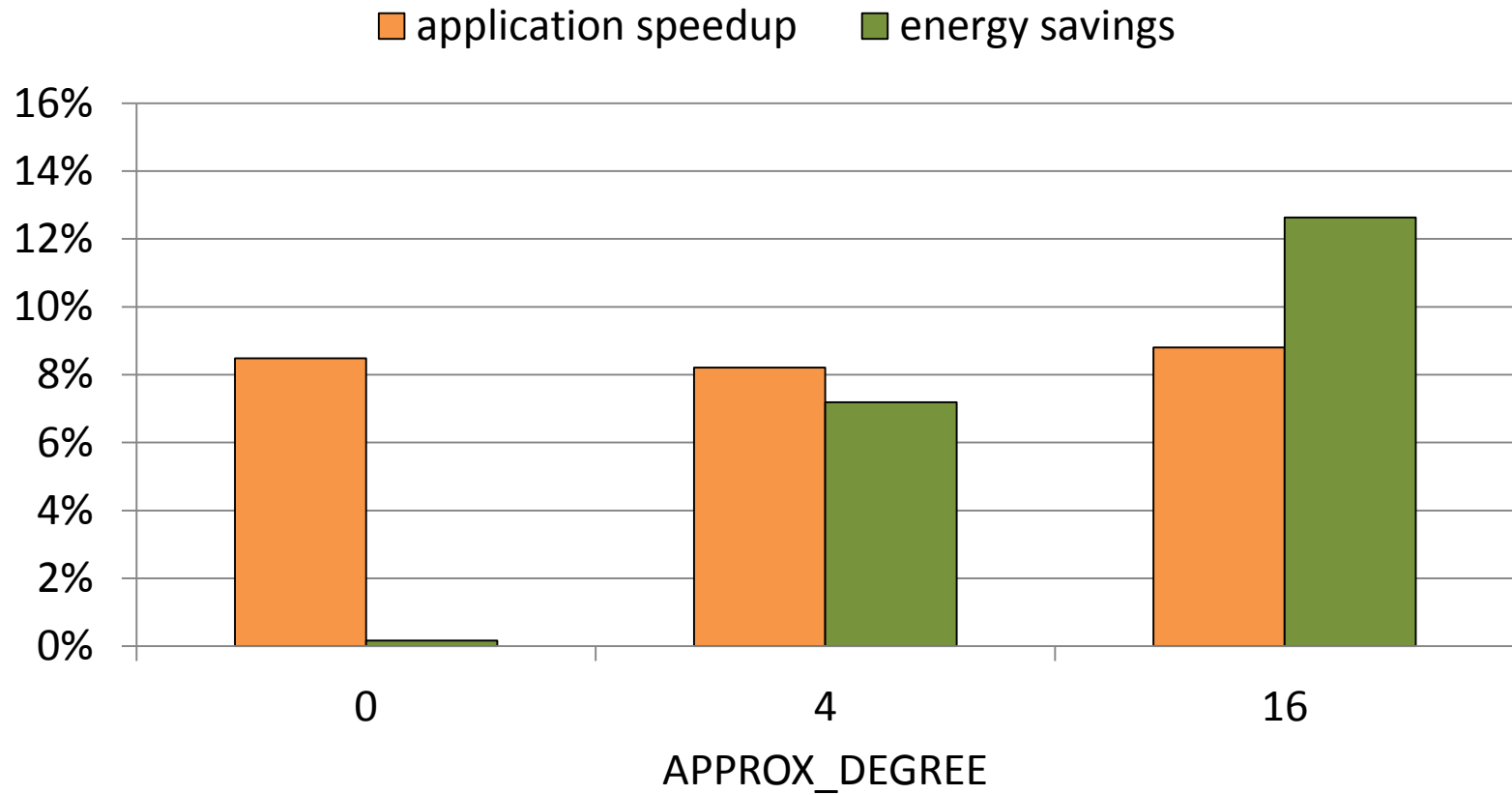
## Full-system simulation

- FeS2 cycle-level x86 simulator [Neelakantam, ASPLOS 2008]

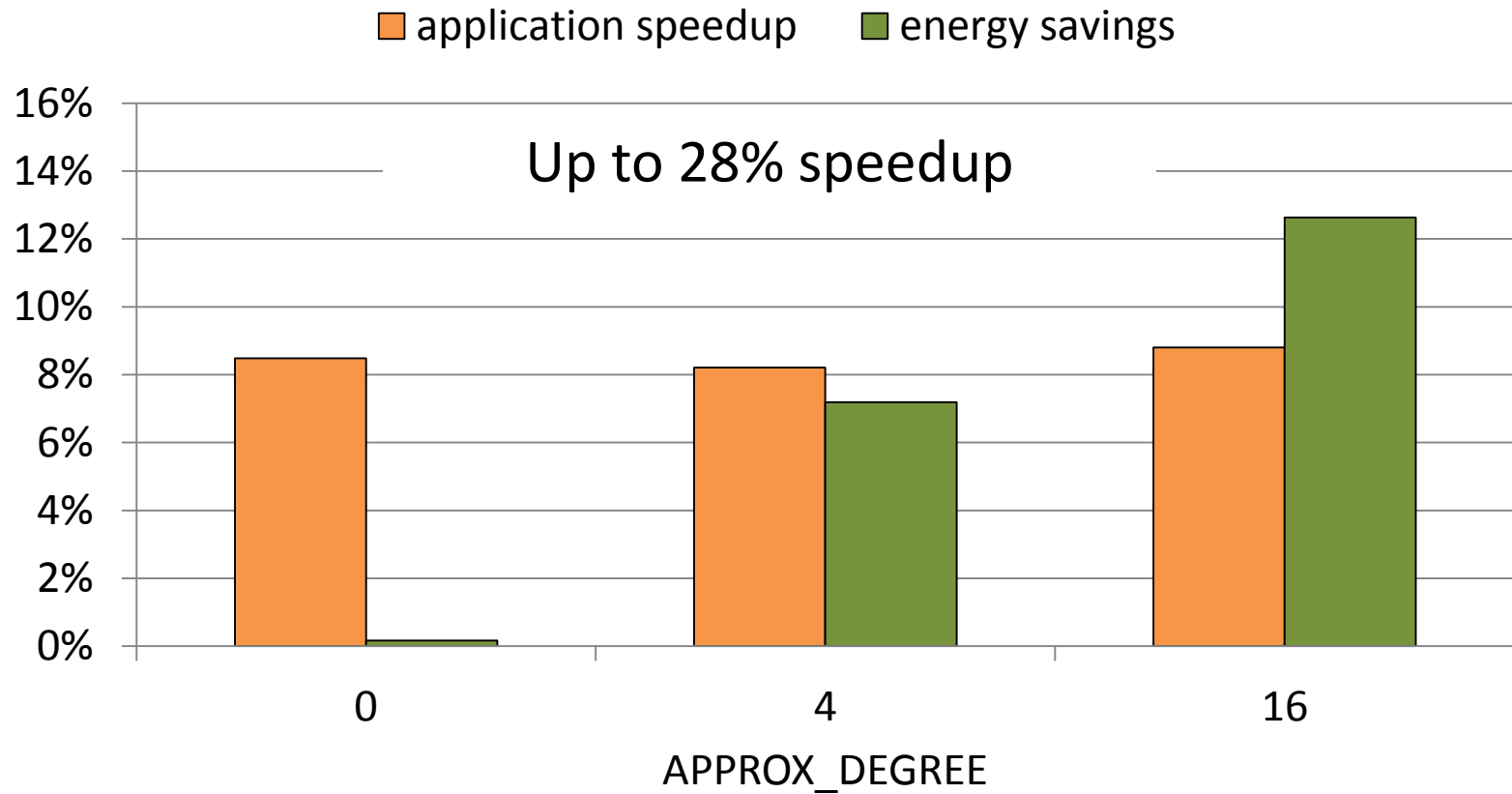
## Approximator, cache and memory energy consumption

- CACTI modeling tool [Thoziyoor, HP 2008]

# Evaluation

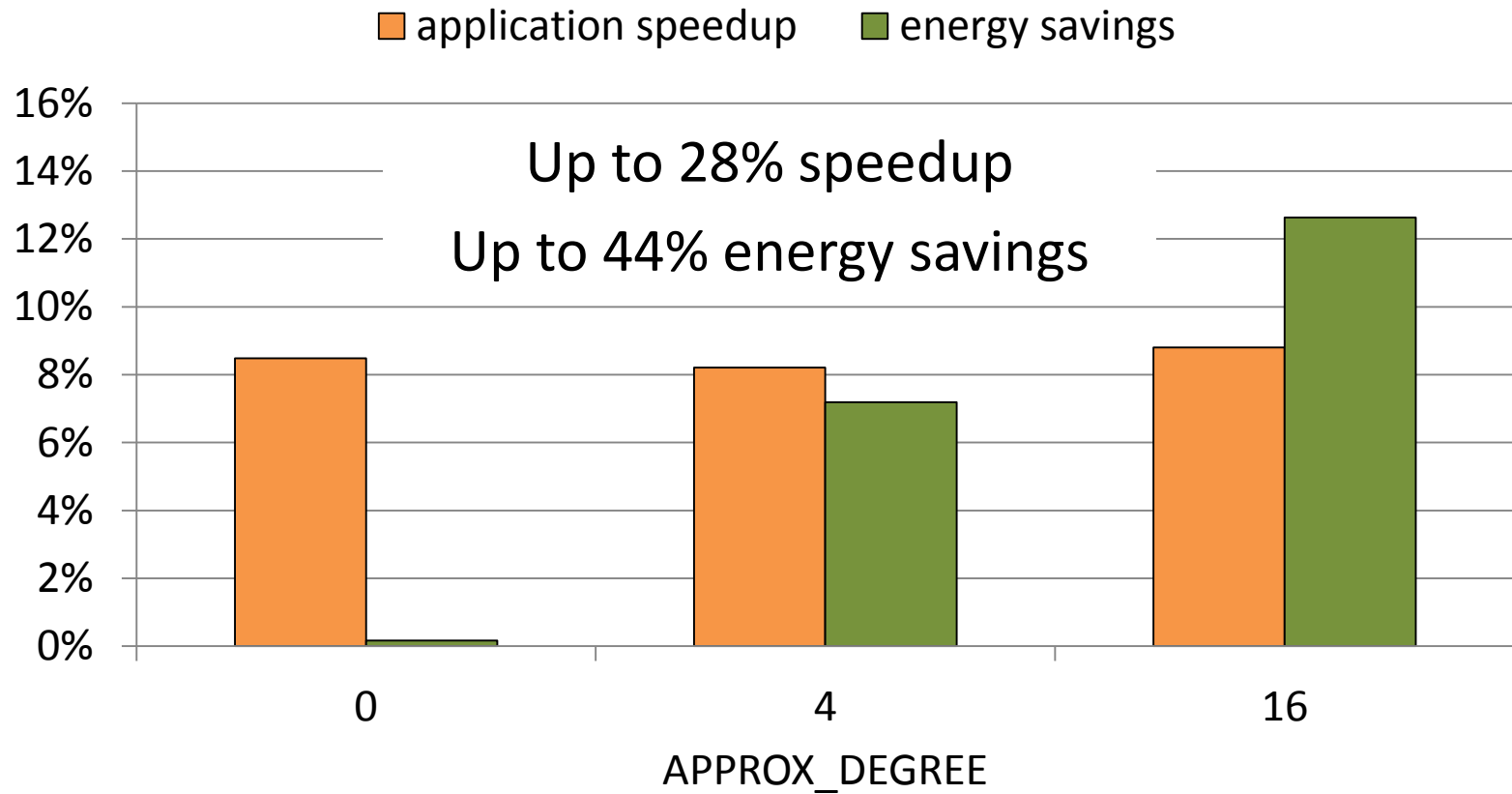


# Evaluation

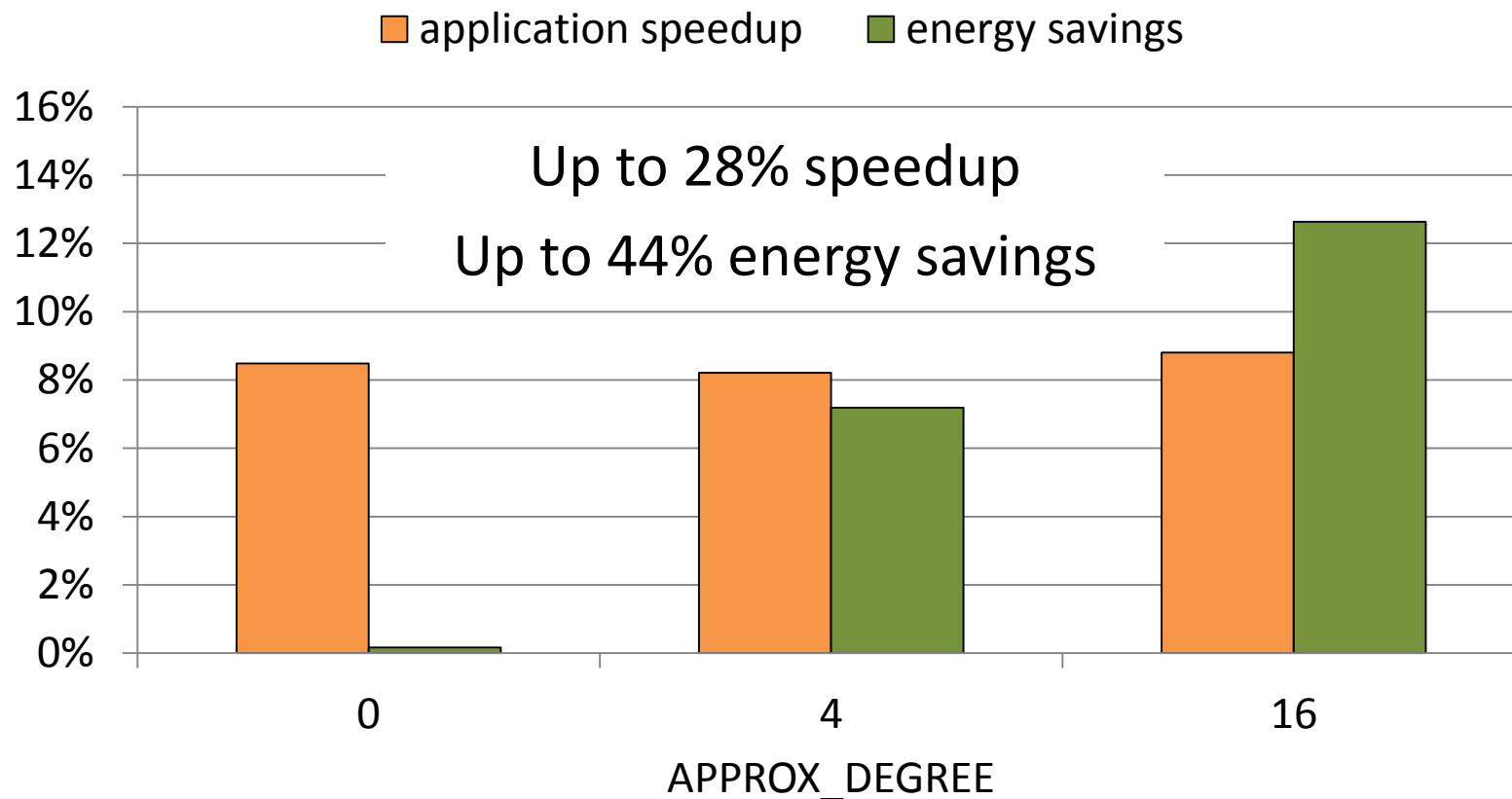




# Evaluation



# Evaluation



Reduces L1-D MPKI by 30% over traditional value predictor and prefetcher.

# Conclusion

## Load Value Approximation

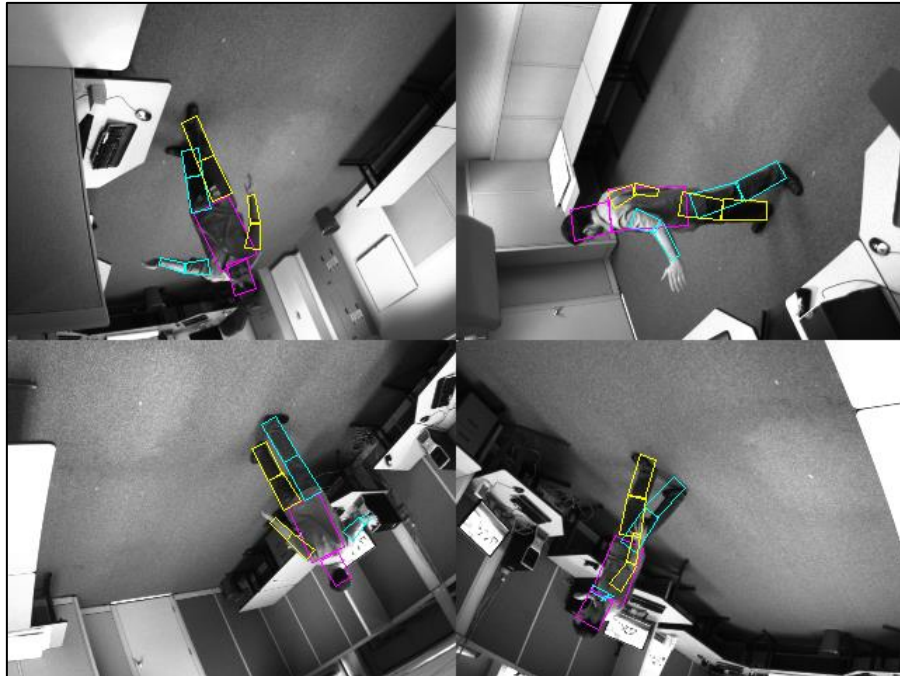
- Approximate Value Locality
- Non-Speculative
- Relaxed Confidence Windows
- Approximation Degree

↑ performance

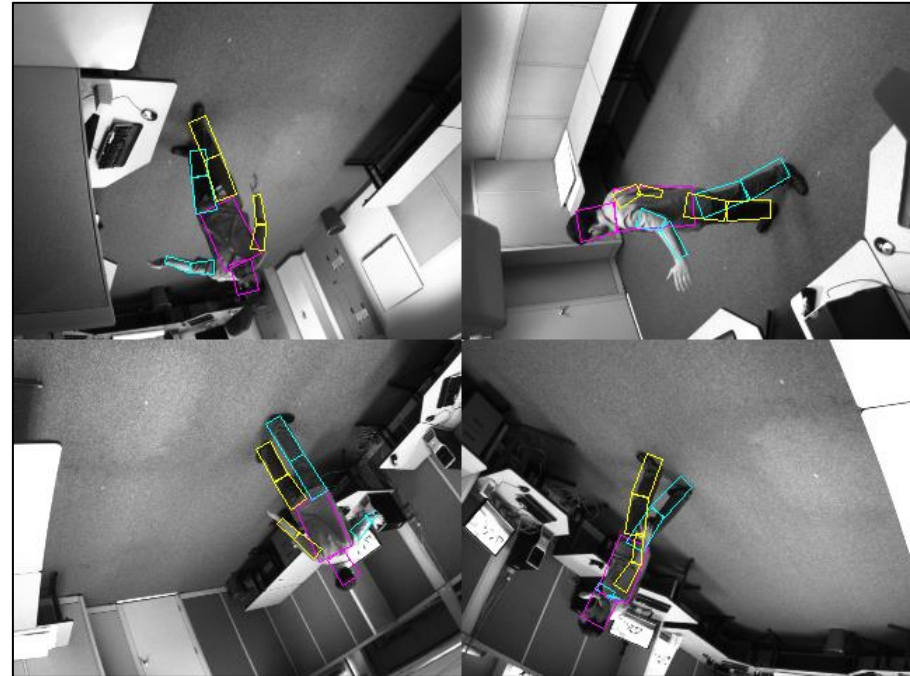
↓ energy consumption

low output error

# Conclusion



baseline (precise)



load value approximation