

Driving Cache Replacement with ML-based LeCaR

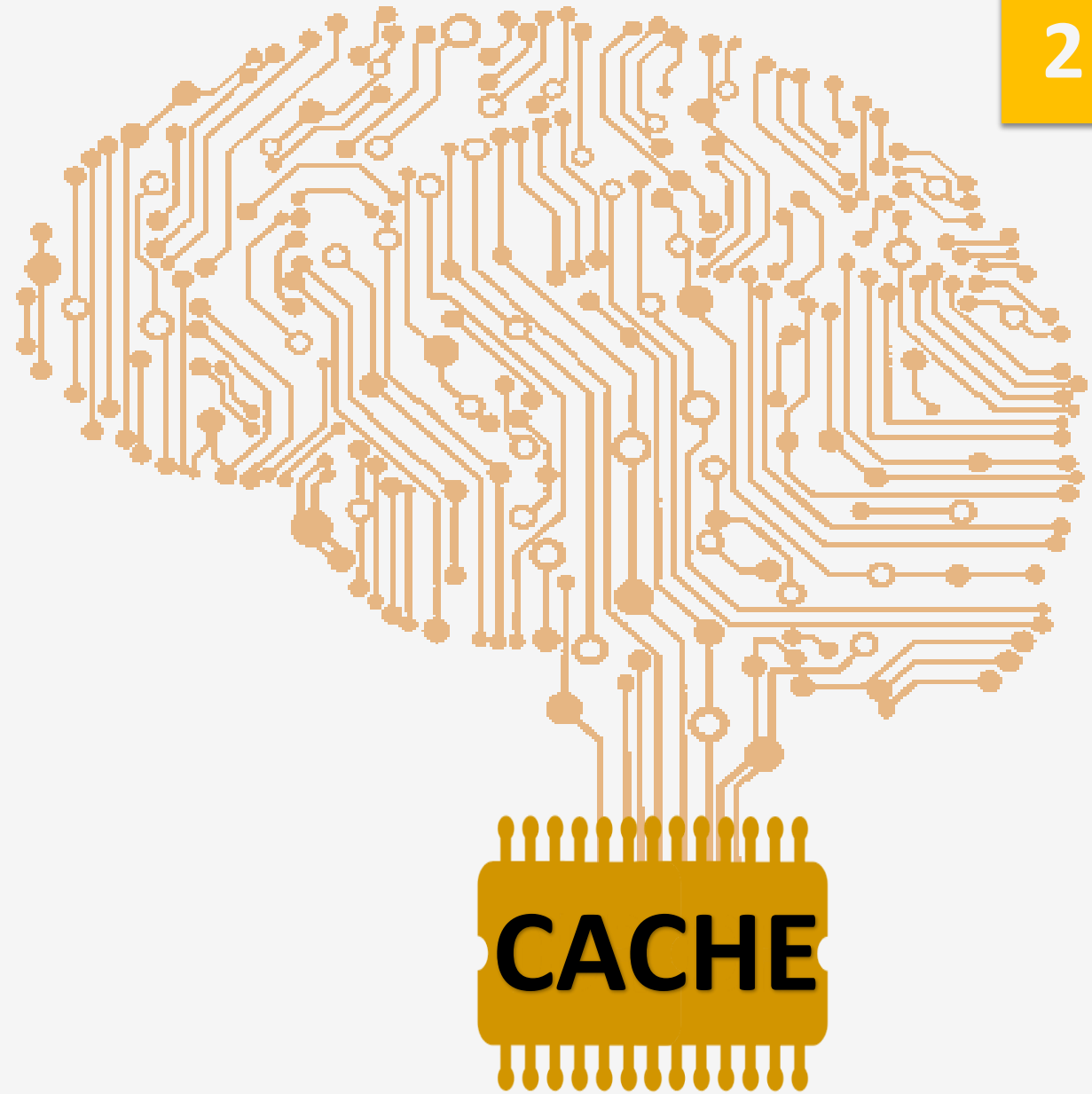
Giuseppe Vietri[†], Liana V. Rodriguez[†], Wendy A. Martinez[†], Steven Lyons[†], Jason Liu[†], Raju Rangaswami[†], Ming Zhao[‡], Giri Narasimhan[†]

[†] Florida International University

[‡] Arizona State University

Reinforcement Learning On Cache Replacement Algorithms

2



Outline

- ⚙ **Previous work**
- ⚙ **Algorithm**
- ⚙ **Results**
- ⚙ **Conclusion**

Hit-rate Performance **vs. ARC**

Cache Size	Non-Parameterized					Adaptive	Fixed-Parameter		
	LRU	LFU	FBR	LIRS	MQ	ARC	LRU-2	2Q	LRFU
1000	-6.1	-10.95	-1.97	-4.13	-1.07	0	0.37	1.55	1.59
2000	-3.61	-10.87	-2.1	-3.57	-1.89	0	-0.26	0.45	0.03
5000	-1.6	-10.49	-1.72	-0.86	-0.86	0	-0.47	0.45	1.48
10000	-1.17	-9.72	0.43	-1.52	-0.79	0	0.55	0.71	1.67
15,000	-0.77	-9.18	0.26	-1.41	-0.59	0	-0.18	0.42	1.66

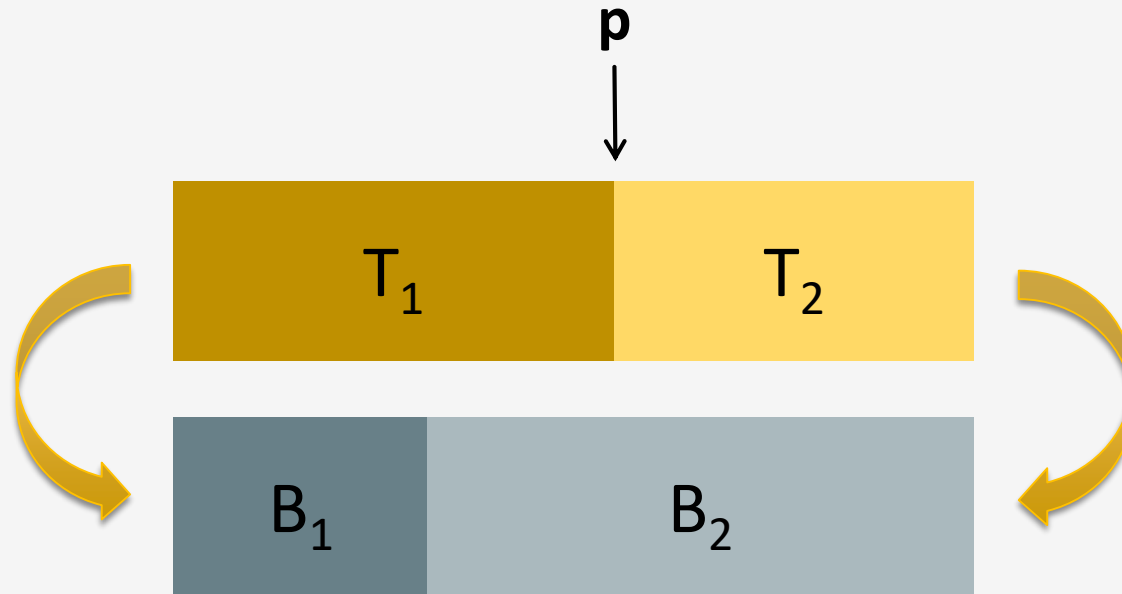


Worse than ARC



Better than ARC

Adaptive Replacement Cache (ARC)

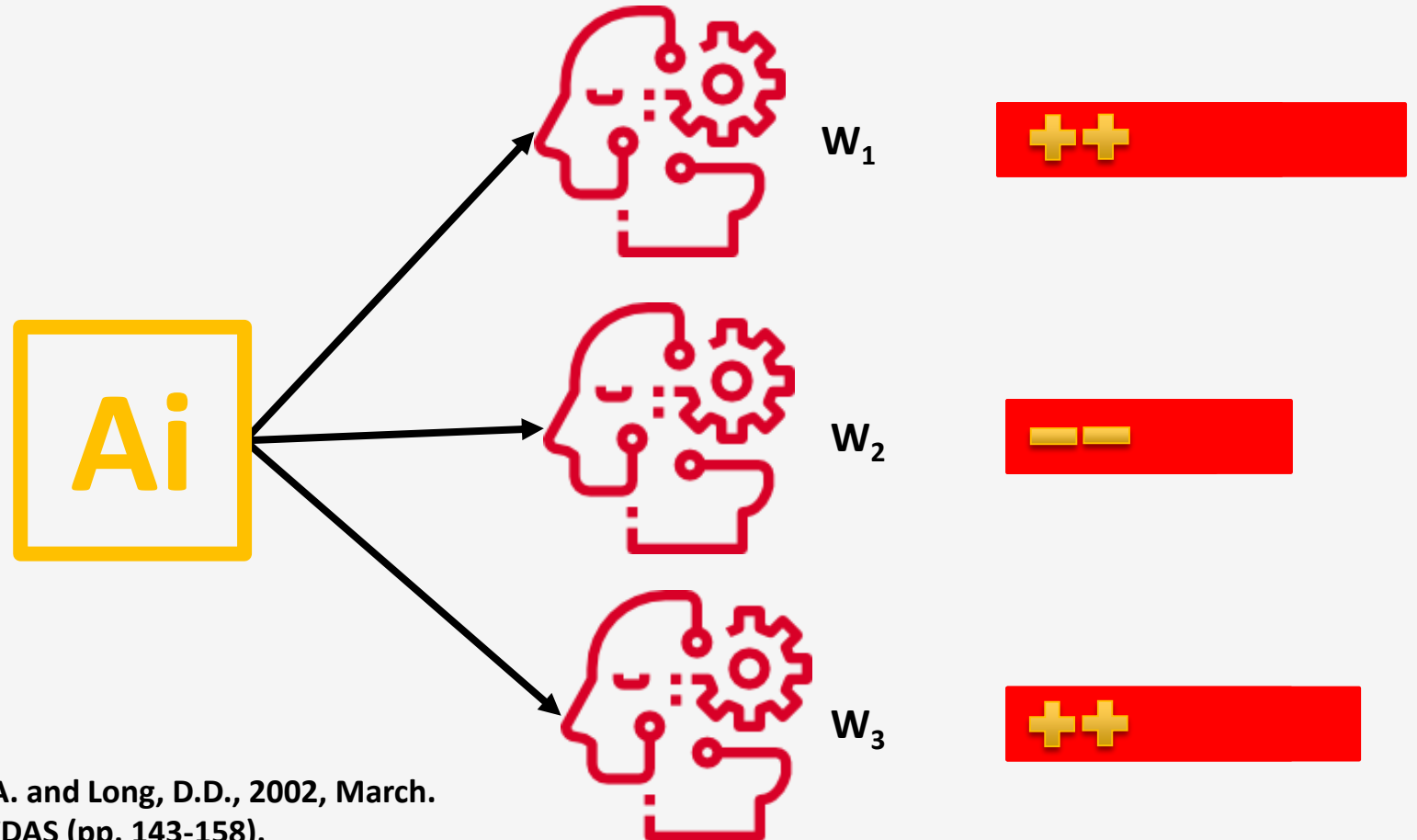


Strengths of ARC

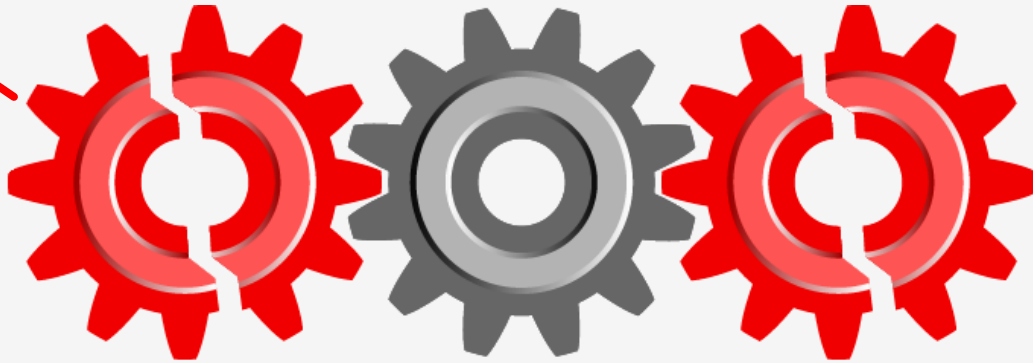
- ⚙️ Manages both **recent** items as well as **frequent** items
- ⚙️ Dynamically **adapts** – Self-tuning
- ⚙️ **Low overhead**
- ⚙️ Competitive Hit-Rate **performance**

Conventional Online Learning Systems

1. Choose Expert
2. Follow Advice
3. Adjust Weights
4. Repeat



Not efficient



Not competitive

**Disadvantages
of current
ML-based
methods**

ML-based LeCaR

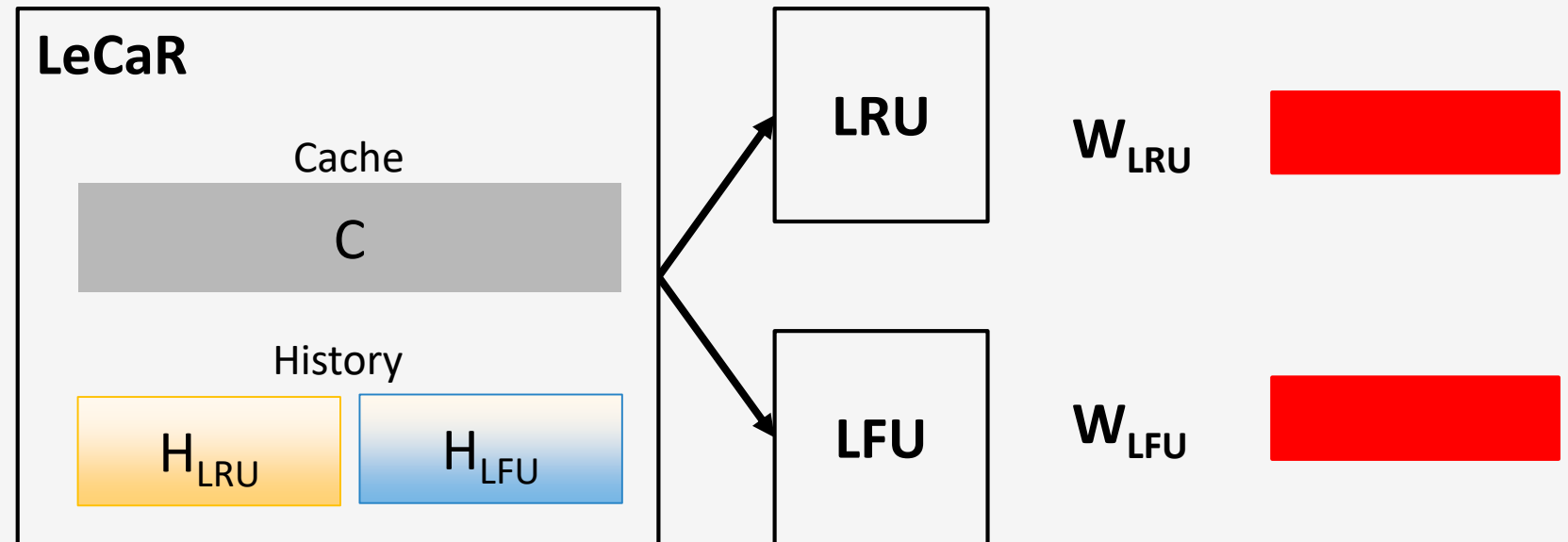
- ⚙ **LRU** and **LFU**
- ⚙ **Efficiency**
- ⚙ **Learning** the optimal mix.


The LeCaR approach

⚙ **LRU** and **LFU**

⚙ **History**

⚙ **Regret**



 **Input:** requested page q
If q in C **then**
 $C.Update(q)$
else
 if q is in H_{LRU} **then**
 $H_{LRU}.Delete(q)$
 if q is in H_{LFU} **then**
 $H_{LFU}.Delete(q)$
 $UpdateWeights(q)$
 if C is full **then**
 action = $(LRU, LFU) \sim (w_{LRU}, w_{LFU})$
 if act == LRU **then**
 if H_{LRU} is full **then**
 $H_{LRU}.Delete(LRU(H_{LRU}))$
 $H_{LRU}.Add(LRU(C))$
 $C.Delete(LRU(C))$
 else
 if H_{LFU} is full **then**
 $H_{LFU}.Delete(LFU(H_{LFU}))$
 $H_{LFU}.Add(LFU(C))$
 $C.Delete(LFU(C))$
 $C.Add(q)$

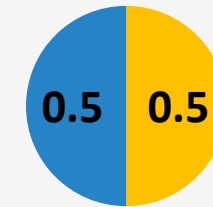
LeCaR (Hit)



11

Request

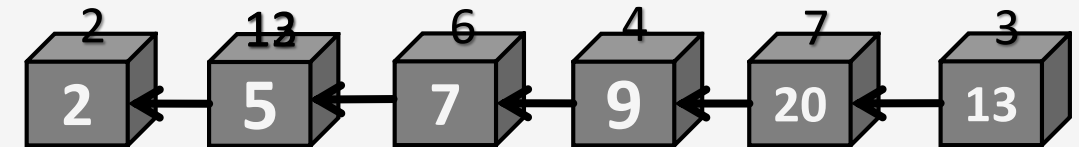


Weights

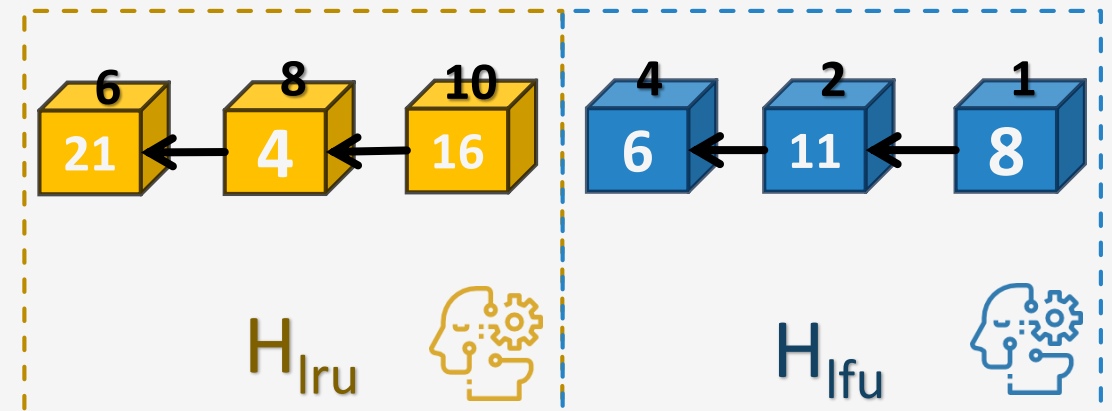


 w_{lru}
 w_{lfu}

Cache



History

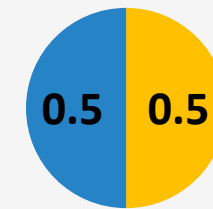


→ **Input:** requested page q
If q in C **then**
 $C.Update(q)$
else
 if q is in H_{LRU} **then**
 $H_{LRU}.Delete(q)$
 if q is in H_{LFU} **then**
 $H_{LFU}.Delete(q)$
 $UpdateWeights(q)$
 if C is full **then**
 action = $(LRU, LFU) \sim (w_{LRU}, w_{LFU})$
 if act == LRU **then**
 if H_{LRU} is full **then**
 $H_{LRU}.Delete(LRU(H_{LRU}))$
 $H_{LRU}.Add(LRU(C))$
 $C.Delete(LRU(C))$
 else
 if H_{LFU} is full **then**
 $H_{LFU}.Delete(LFU(H_{LFU}))$
 $H_{LFU}.Add(LFU(C))$
 $C.Delete(LFU(C))$
 $C.Add(q)$

LeCaR (**Miss not in History**)

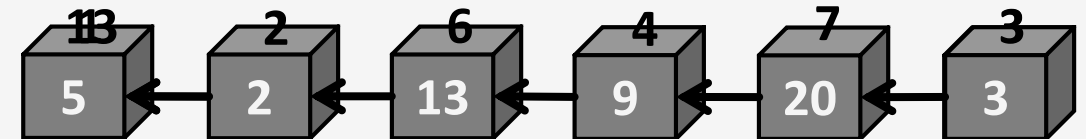
12

Request

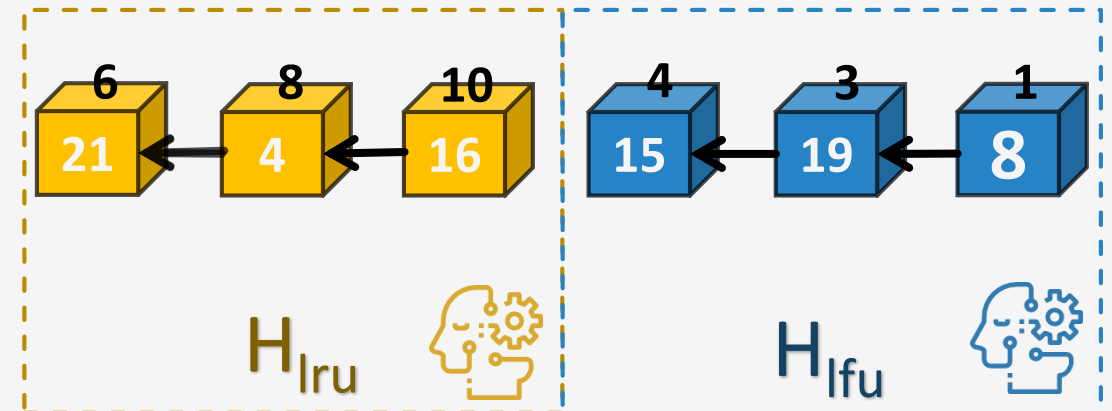


■ w_{lru}
■ w_{lfu}

Cache



History



Weight Update

α = Learning rate

	LRU Regret	LFU Regret
Update	$w_{\text{LRU}} := \alpha * w_{\text{LRU}}$	$w_{\text{LFU}} := \alpha * w_{\text{LFU}}$

→ **Input:** requested page q

```

if  $q$  in  $C$  then
   $C.Update(q)$ 
else
  if  $q$  is in  $H_{LRU}$  then
     $H_{LRU}.Delete(q)$ 
  if  $q$  is in  $H_{LFU}$  then
     $H_{LFU}.Delete(q)$ 
   $UpdateWeights(q)$ 
  if  $C$  is full then
    action =  $(LRU, LFU) \sim (w_{LRU}, w_{LFU})$ 
    if act == LRU then
      if  $H_{LRU}$  is full then
         $H_{LRU}.Delete(LRU(H_{LRU}))$ 
       $H_{LRU}.Add(LRU(C))$ 
       $C.Delete(LRU(C))$ 
    else
      if  $H_{LFU}$  is full then
         $H_{LFU}.Delete(LFU(H_{LFU}))$ 
       $H_{LFU}.Add(LFU(C))$ 
       $C.Delete(LFU(C))$ 
   $C.Add(q)$ 

```

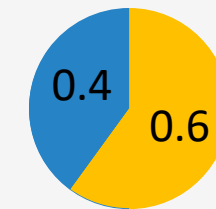
LeCaR (**Miss in History**)

14

Request

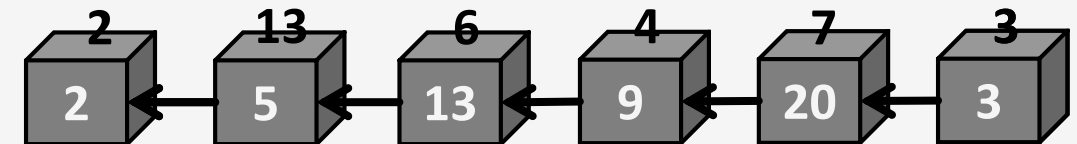


Weights

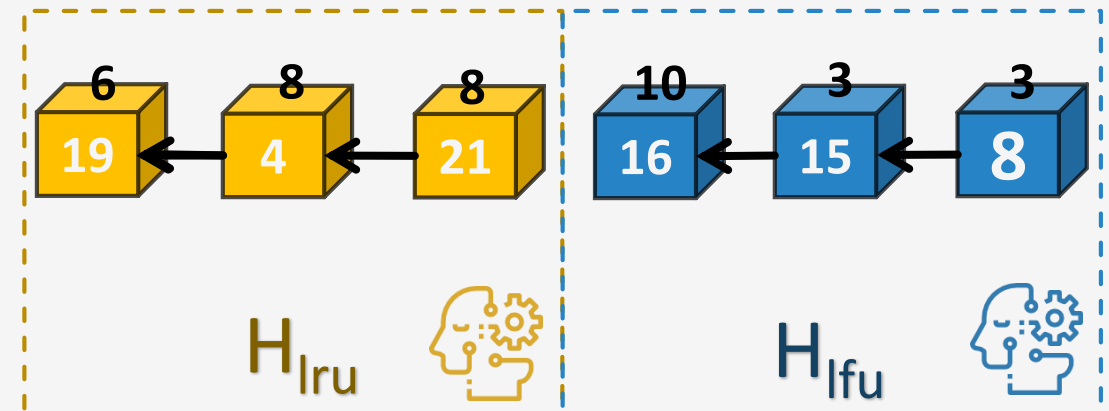


■ W_{lru}
■ W_{lfu}

Cache



History



Experiments

- ⚙️ **8** Workloads
- ⚙️ 3 days of data[FIU datasets]
- ⚙️ Small, Medium and Large **Cache Sizes**
- ⚙️ Fixed **Learning Rate**

Data used for the Experiments

**Casa, Ikki,
Madmax
and Topgun**

Four different end-user/developer home directories

Online

Departments online course management system

Webresearch

Document store for research projects

Webmail

Mail server of the FIU Computer and Engineering department using Postfix

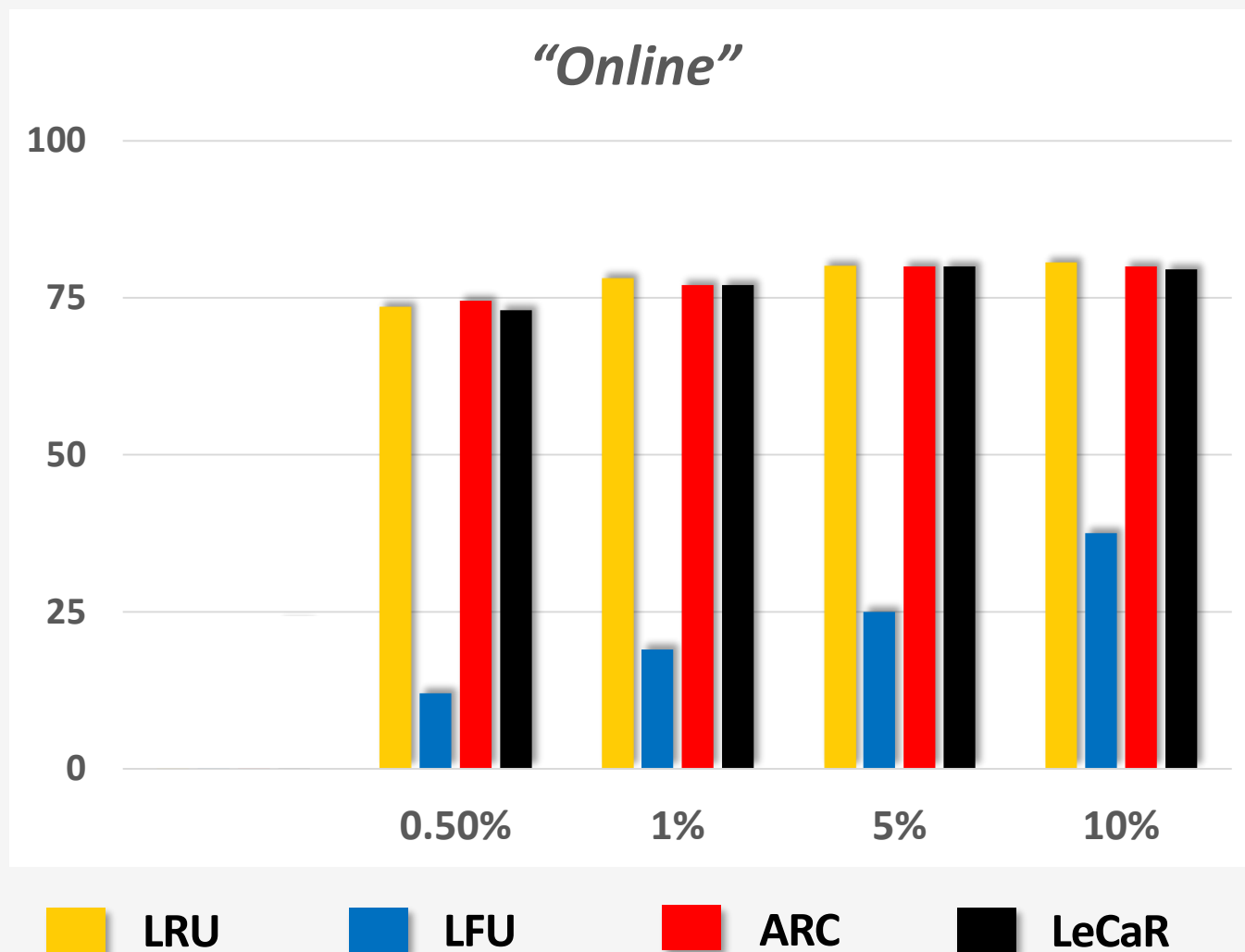
Webuser

Web server hosting faculty, staff, and graduate student web sites



*Collected at the School of
Computing and Information
Sciences at FIU.*

Results



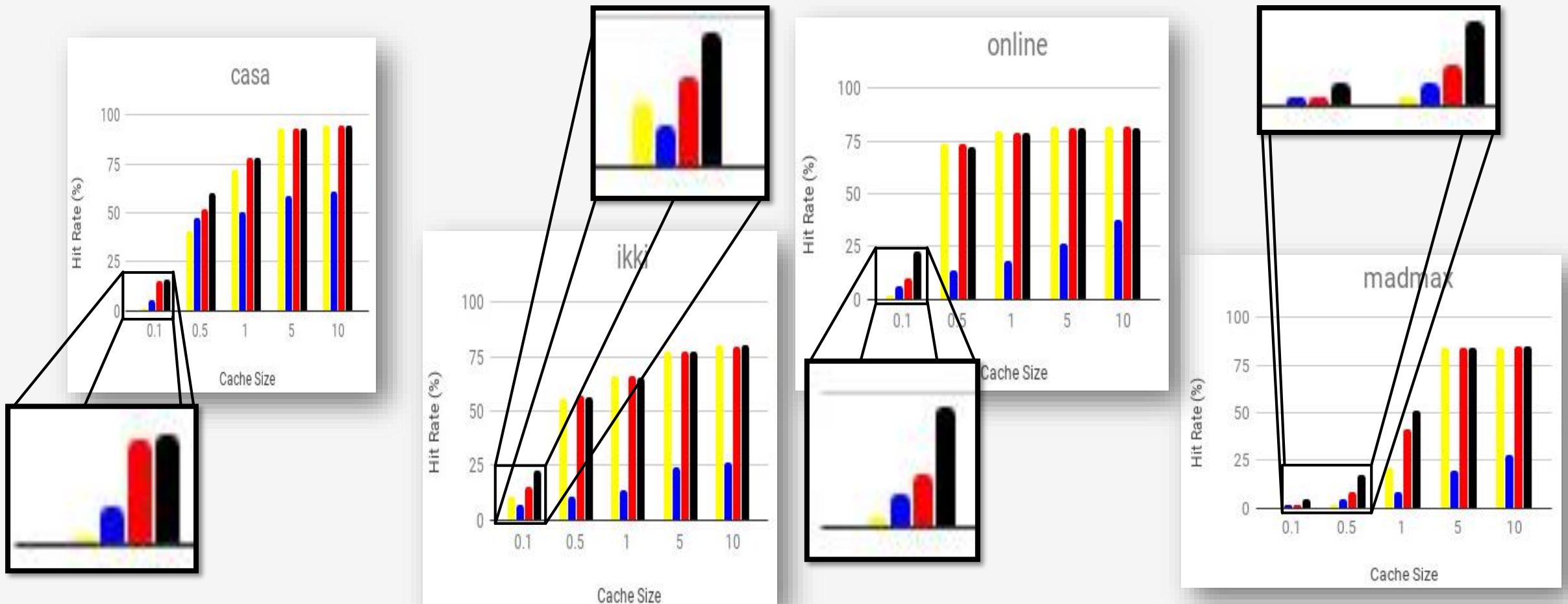
Results

LRU

LFU

ARC

LeCaR



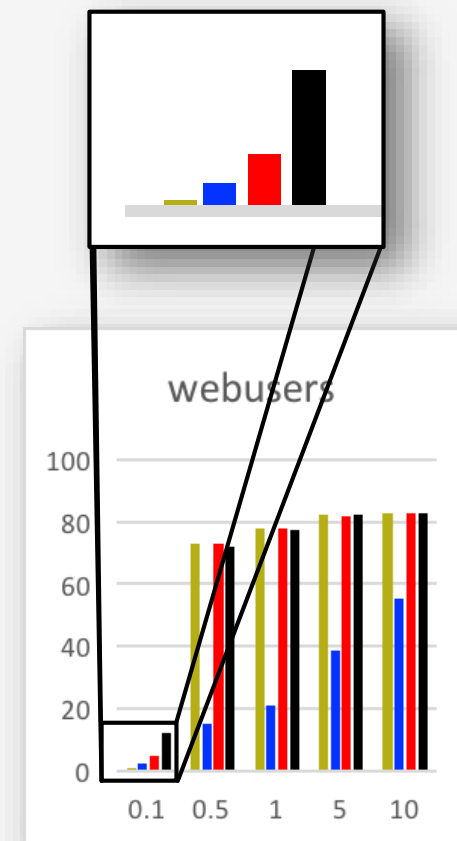
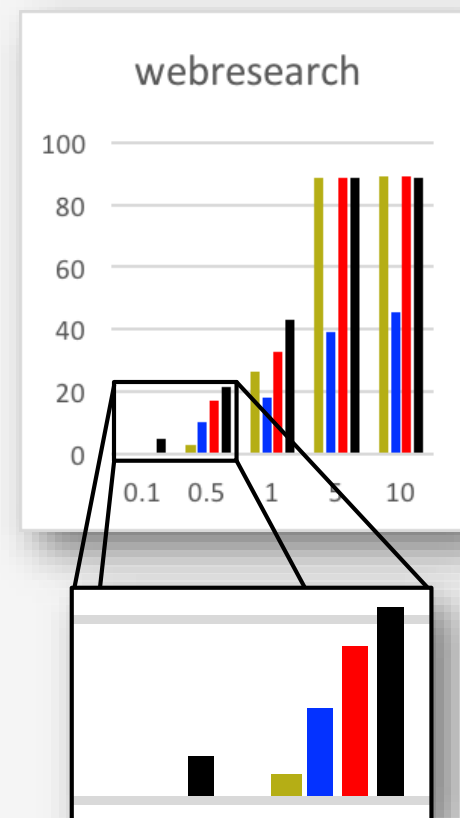
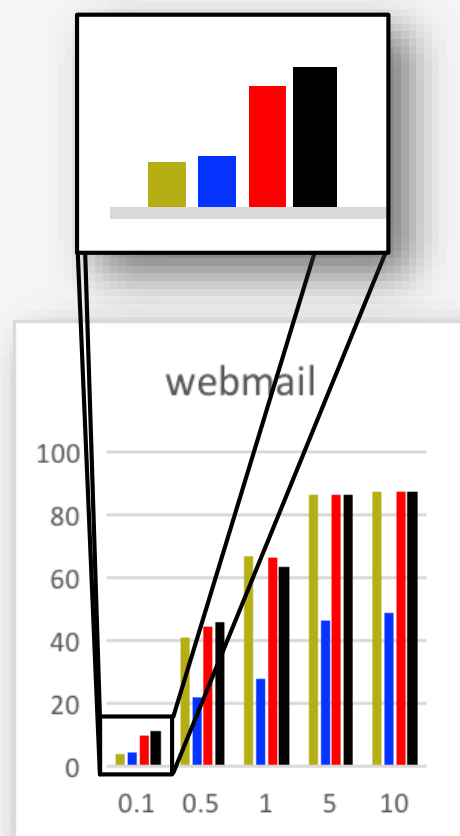
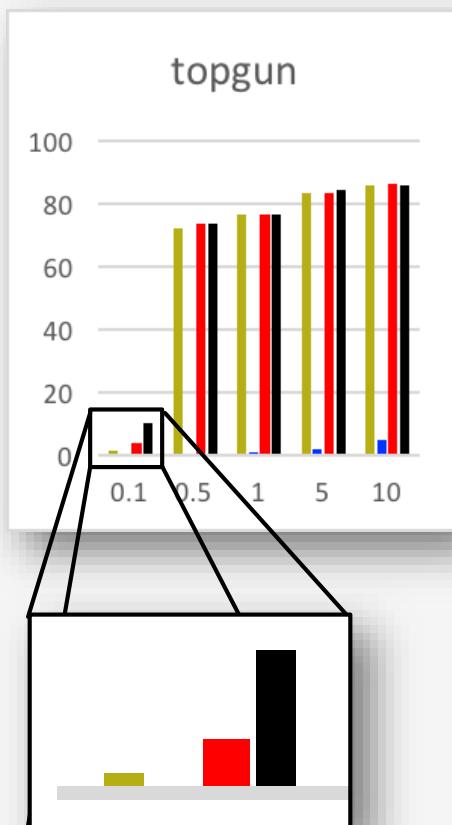
Results

LRU

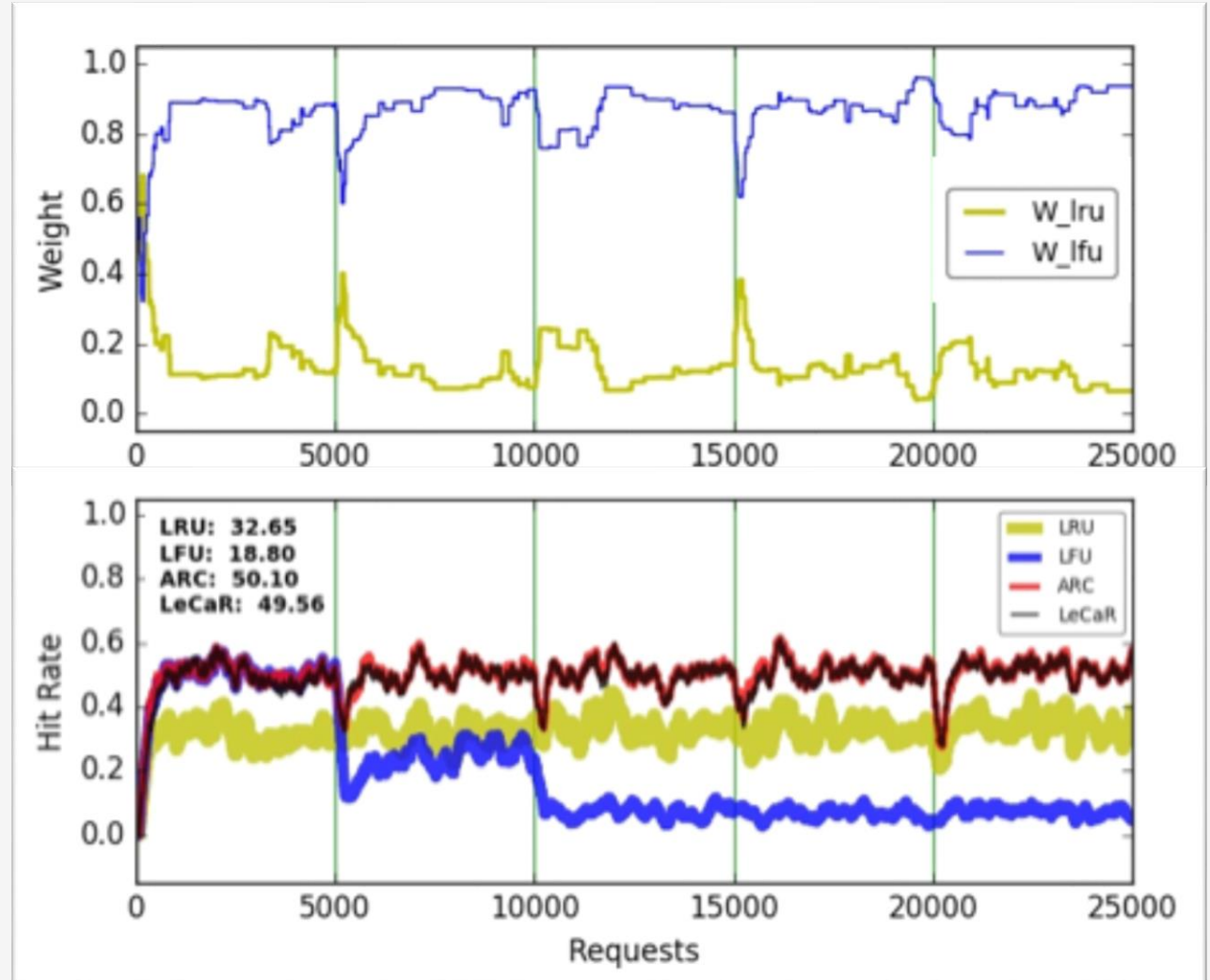
LFU

ARC

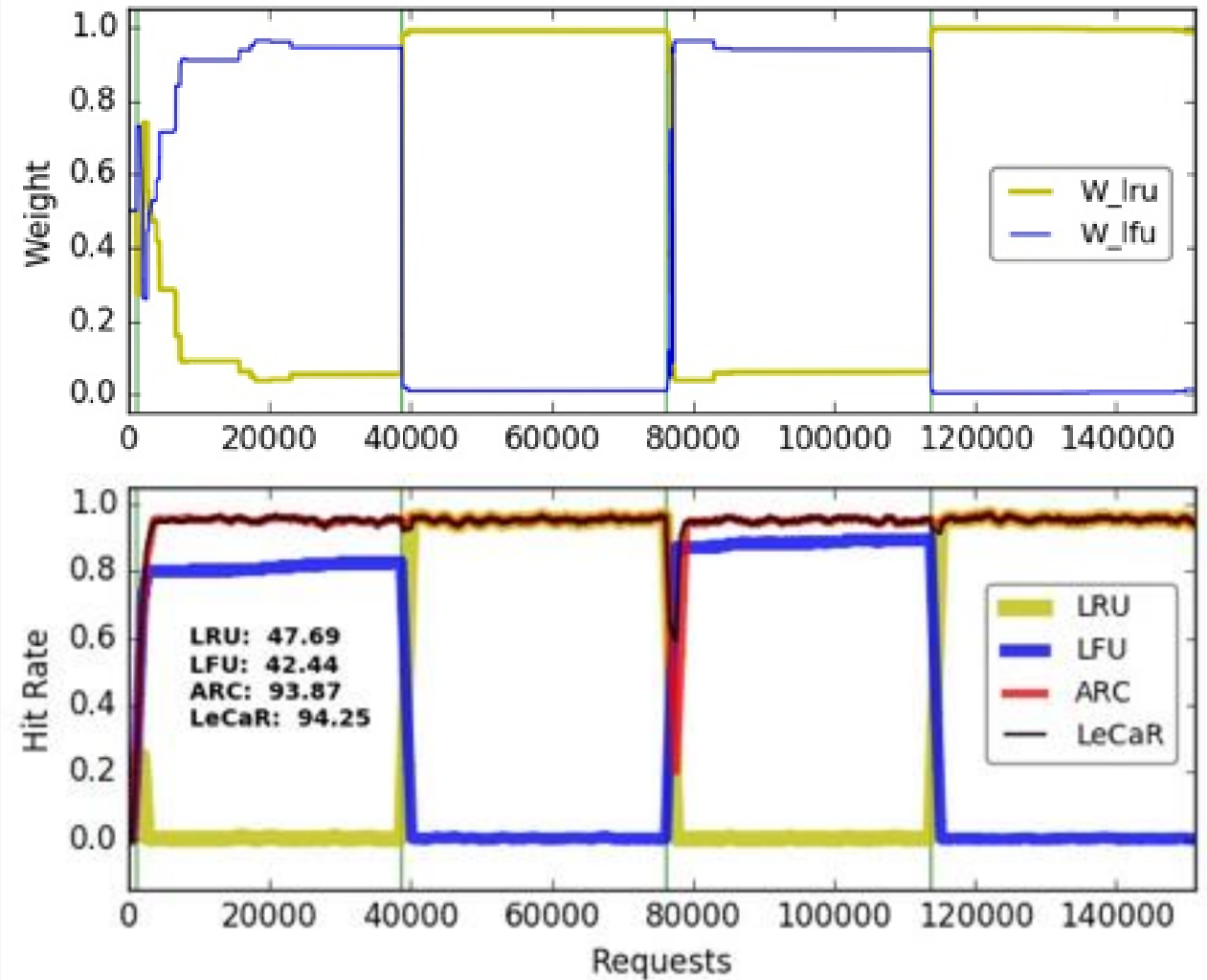
LeCaR



Synthetic Data



Synthetic Data

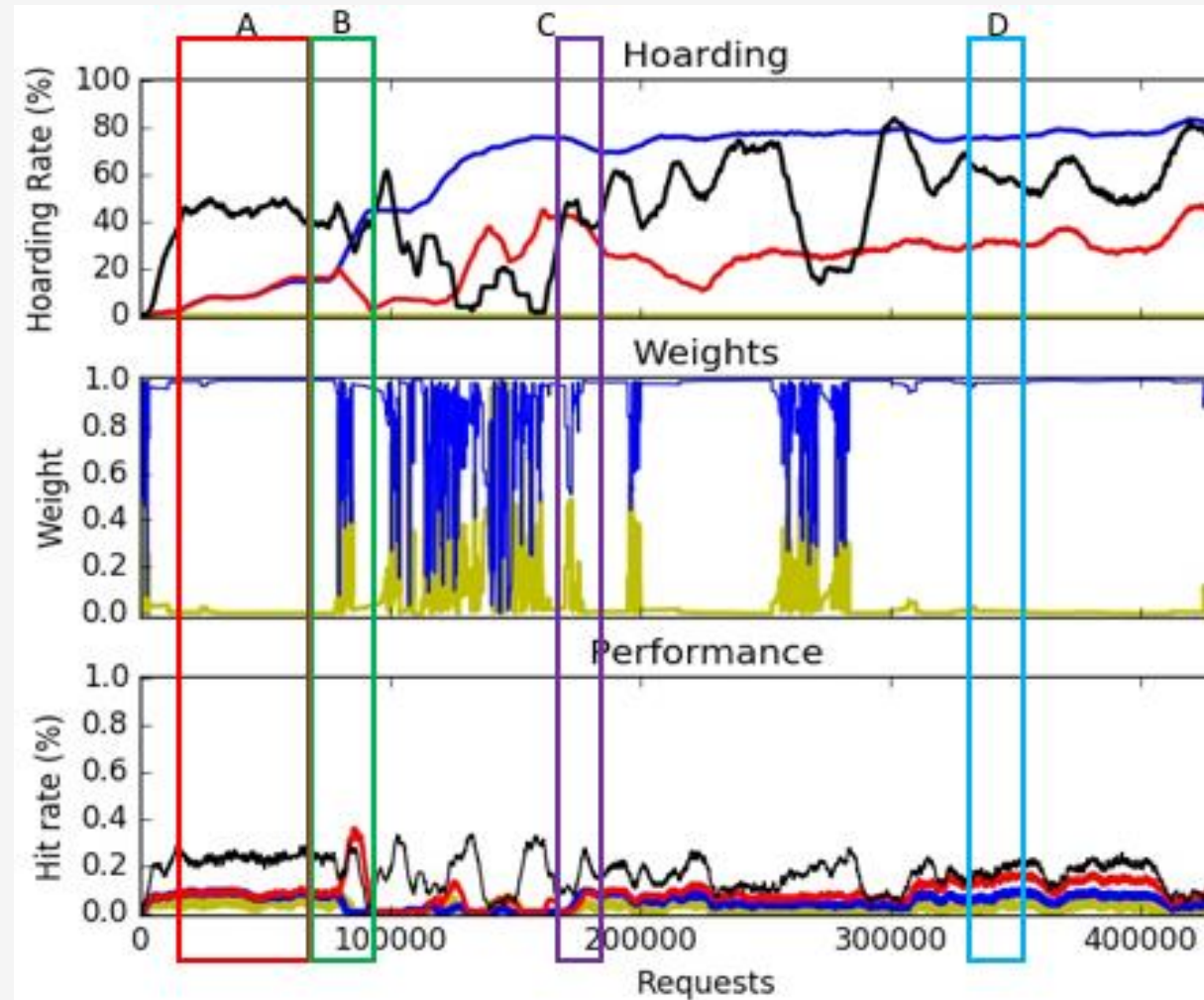


Hoarding Rate

- ⚙ Definition
- ⚙ A: Stable Period
- ⚙ B: LFU gets penalized

Hoarded Page:

- ⚙ Accessed at least twice
- ⚙ Not among the last $2N$ unique pages



Conclusion

- ⚙ ML-based **LeCaR**
- ⚙ Small cache size
- ⚙ **Real + Synthetic** experiments