# Design and Analysis of Real-Time Systems
# Caches in WCET Analysis

Jan Reineke

Department of Computer Science
Saarland University
Saarbrücken, Germany

Advanced Lecture, Summer 2013

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Caches

- How they work:
  - dynamically
  - managed by replacement policy



| | [*ab*] | |
|---|---|---|
| CPU | Cache | Main Memory |

| | | |
|---|---|---|
| Capacity: | 32 KB | 2 MB |
| Latency: | 3 cycles | 100 cycles |

- Why they work: *principle of locality*
  - spatial
  - temporal

# Caches



- How they work:
  - dynamically
  - managed by replacement policy



- Why they work: *principle of locality*
  - spatial
  - temporal

# Caches

- How they work:
  - dynamically
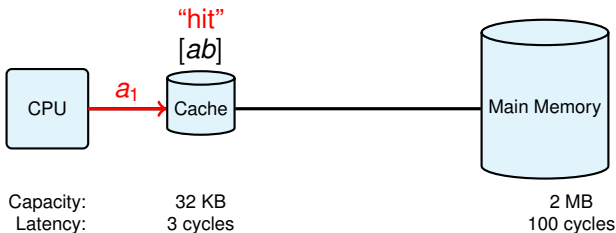  - managed by replacement policy



- Why they work: *principle of locality*
  - spatial
  - temporal

# Caches

- How they work:
  - dynamically
  - managed by replacement policy



- Why they work: *principle of locality*
  - spatial
  - temporal

# Caches

- How they work:
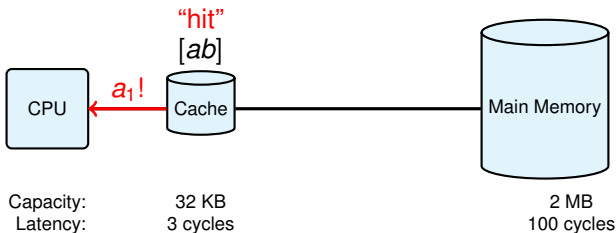  - dynamically
  - managed by replacement policy



- Why they work: *principle of locality*
  - spatial
  - temporal

# Caches

- How they work:
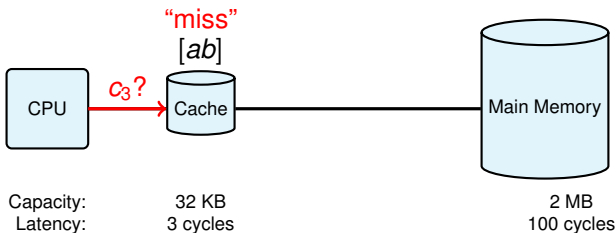  - dynamically
  - managed by replacement policy



- Why they work: *principle of locality*
  - spatial
  - temporal

# Caches

- How they work:
  - dynamically
  - managed by replacement policy

"miss"
[*ac*]



| | | | |
|---|---|---|---|
| Capacity: | 32 KB | | 2 MB |
| Latency: | 3 cycles | | 100 cycles |

- Why they work: *principle of locality*
  - spatial
  - temporal

# Caches

- How they work:
  - dynamically
  - managed by replacement policy



- Why they work: *principle of locality*
  - spatial
  - temporal

# Caches

- How they work:
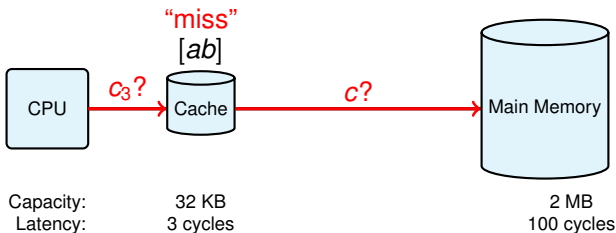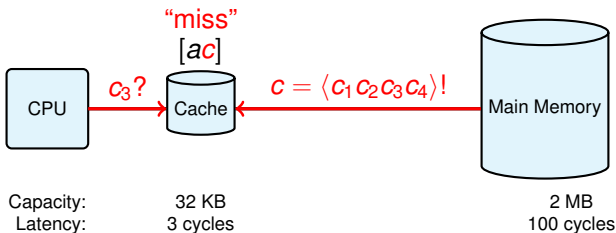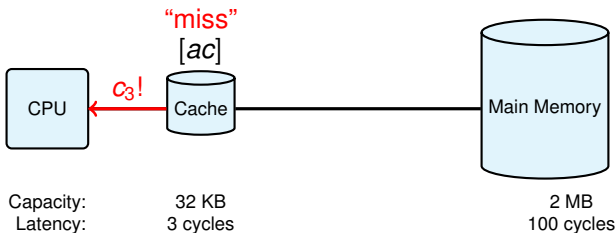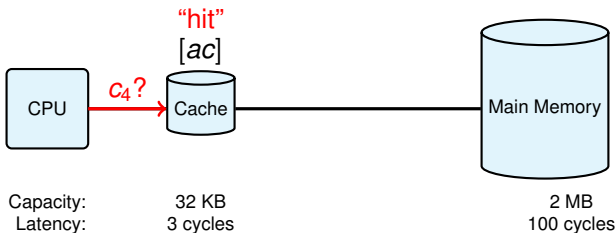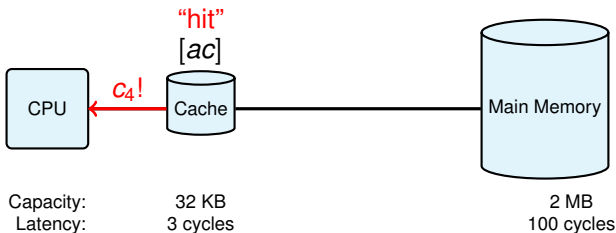  - dynamically
  - managed by replacement policy



| | Capacity: | 32 KB | | 2 MB |
| Latency: | 3 cycles | | 100 cycles |

- Why they work: *principle of locality*
  - spatial
  - temporal

# Fully-Associative Caches

# Set-Associative Caches

*Address:*

$\leftarrow log_2(s) \rightarrow$ $\leftarrow log_2(8 * b) \rightarrow$

| Tag | Index | Block offset |

*Cache Set:*

| Tag | Data Block |
| Tag | Data Block |
... 
| Tag | Data Block |

...

*Cache Set:*

| Tag | Data Block |
| Tag | Data Block |
...
| Tag | Data Block |

$k$

$s$

=?  Yes: Hit!

No: Miss!

MUX

Data

Special cases:

- direct-mapped cache: only one line per cache set
- fully-associative cache: only one cache set

# Cache Replacement Policies

- Least-Recently-Used (LRU) used in
    INTEL PENTIUM I and MIPS 24K/34K
- First-In First-Out (FIFO or Round-Robin) used in
    MOTOROLA POWERPC 56X, INTEL XSCALE, ARM9, ARM11
- Pseudo-LRU (PLRU) used in
    INTEL PENTIUM II-IV and POWERPC 75X
- Most Recently Used (MRU) as described in literature

Each cache set is treated independently:
$\longrightarrow$ Set-associative caches are compositions of fully-associative caches.

のため

# Cache Analysis

Two types of cache analyses:

1. Local guarantees: classification of individual accesses
   - May-Analysis $\longrightarrow$ Overapproximates cache contents
   - Must-Analysis $\longrightarrow$ Underapproximates cache contents
2. Global guarantees: bounds on cache hits/misses

- Cache analyses almost exclusively for LRU
- In practice: FIFO, PLRU, . . .

Always a cache hit/always a miss?

SAARLAND
UNIVERSITY
COMPUTER SCIENCE



Always a cache hit/always a miss?

1. Initial cache contents unknown.

2. Different paths lead to these points.

3. Cannot resolve address of $z$.

SAARLAND
UNIVERSITY
COMPUTER SCIENCE



*Collecting Semantics* =
set of states at each program point that
any execution may encounter there

Two approximations:

| | Collecting Semantics | uncomputable |
|---|---|---|
| $\subseteq$ | Cache Semantics | computable |
| $\subseteq$ | $\gamma$(Abstract Cache Sem.) | efficiently computable |

*Collecting Semantics* =
set of states at each program point that
any execution may encounter there

Two approximations:

$\phantom{\subseteq}$ Collecting Semantics     uncomputable

$\subseteq$ Cache Semantics     computable

$\subseteq$ $\gamma$(Abstract Cache Sem.) efficiently
                            computable

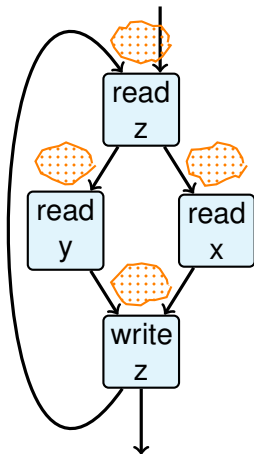# Deriving Invariants about Cache States using Abstract Interpretation



*Collecting Semantics =*
set of states at each program point that any execution may encounter there

Two approximations:

| | Collecting Semantics | uncomputable |
| --- | --- | --- |
| $\subseteq$ | Cache Semantics | computable |
| $\subseteq$ | $\gamma$(Abstract Cache Sem.) | efficiently computable |

# Deriving Invariants about Cache States using Abstract Interpretation



*Collecting Semantics =*
set of states at each program point that any execution may encounter there

Two approximations:

|  | Collecting Semantics | uncomputable |
|---|---|---|
| $\subseteq$ | Cache Semantics | computable |
| $\subseteq$ | $\gamma$(Abstract Cache Sem.) | efficiently computable |

# Deriving Invariants about Cache States using Abstract Interpretation



*Collecting Semantics =*
set of states at each program point that
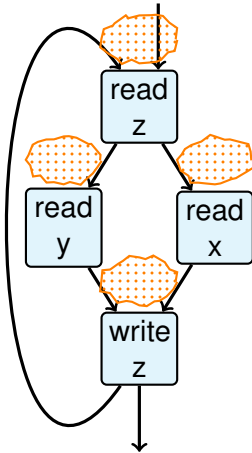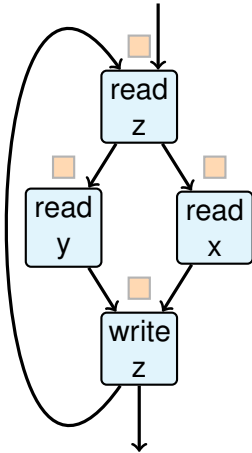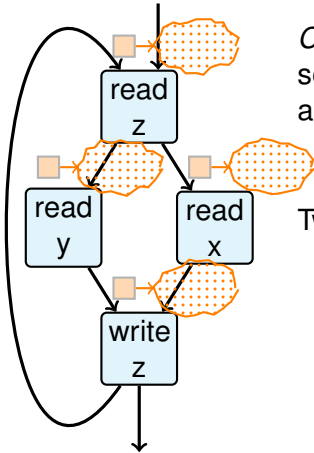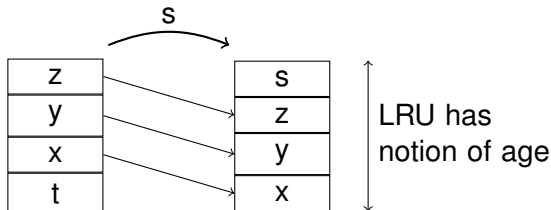any execution may encounter there

Two approximations:

| Collecting Semantics | uncomputable |
| $\subseteq$ Cache Semantics | computable |
| $\subseteq$ $\gamma$(Abstract Cache Sem.) | efficiently computable |

# Least-Recently-Used (LRU): Concrete Behavior

- Used to predict *cache hits*.
- Maintains *upper bounds on ages* of memory blocks.
- Upper bound $\leq$ associativity $\longrightarrow$ memory block definitely cached.

## Example

Abstract state:

. . . and its interpretation:

Describes the set of all concrete cache states in which $x$, $s$, and $t$ occur,

| | |
|---|---|
| {x} | age 0 |
| {} | |
| {s,t} | |
| {} | age 3 |

- $x$ with an age of 0,
- $s$ and $t$ with an age not older than 2.

$$\gamma([\{x\}, \{\}, \{s, t\}, \{\}]) = \{[x, s, t, a], [x, t, s, a], [x, s, t, b], \ldots\}$$

# Sound Update – Local Consistency



concrete cache states

concrete cache states

"Potential Cache Miss":

z

| {x} | → | {z} |
| {} | | {x} |
| {s,t} | | {} |
| {} | | {s,t} |

"Definite Cache Hit":

s

| {x} | → | {s} |
| {} | | {x} |
| {s,t} | | {t} |
| {} | | {} |

Why does *t* not age in the second case?

# LRU: Must-Analysis: Join

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
- $\gamma(B) \subseteq \gamma(A \sqcup B)$

| {a} |
|-----|
| {} |
| {c,f} |
| {d} |

$\sqcup$

| {c} |
|-----|
| {e} |
| {a} |
| {d} |

| {} |
|-----|
| {} |
| {a,c} |
| {d} |

"Intersection + Maximal Age"

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
- $\gamma(B) \subseteq \gamma(A \sqcup B)$

| {a} |
|-----|
| {} |
| {c,f} |
| {d} |

$\sqcup$

| {c} |
|-----|
| {e} |
| {a} |
| {d} |

| {} |
|-----|
| {} |
| {a,c} |
| {d} |

"Intersection + Maximal Age"

# LRU: Must-Analysis: Join

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
- $\gamma(B) \subseteq \gamma(A \sqcup B)$



| {a} |
|-----|
| {} |
| {c,f} |
| {d} |

$\sqcup$

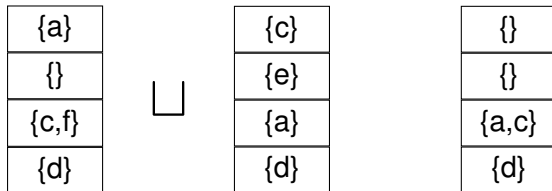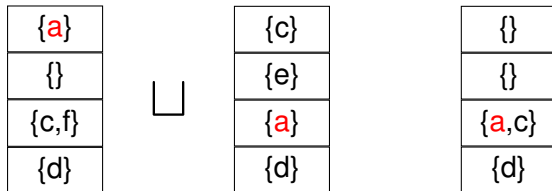| {c} |
|-----|
| {e} |
| {a} |
| {d} |

| {} |
|-----|
| {} |
| {a,c} |
| {d} |

"Intersection + Maximal Age"

# LRU: Must-Analysis: Join

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
- $\gamma(B) \subseteq \gamma(A \sqcup B)$

| {a} |
|-----|
| {} |
| {c,f} |
| {d} |

$\sqcup$

| {c} |
|-----|
| {e} |
| {a} |
| {d} |

| {} |
|-----|
| {} |
| {a,c} |
| {d} |

"Intersection + Maximal Age"

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
- $\gamma(B) \subseteq \gamma(A \sqcup B)$



| {a} |
| {} |
| {c,f} |
| {d} |

$\sqcup$

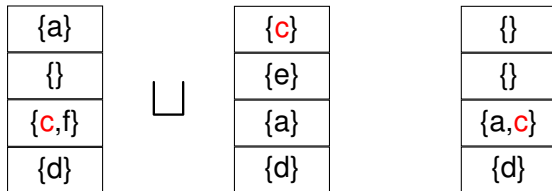| {c} |
| {e} |
| {a} |
| {d} |

| {} |
| {} |
| {a,c} |
| {d} |

"Intersection + Maximal Age"

# LRU: Must-Analysis: Join

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
- $\gamma(B) \subseteq \gamma(A \sqcup B)$



| {a} |
|-----|
| {} |
| {c,f} |
| {d} |

$\sqcup$

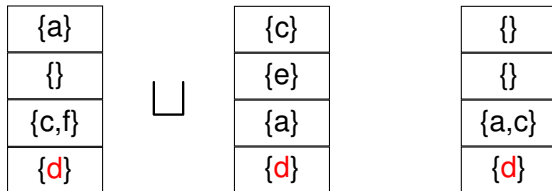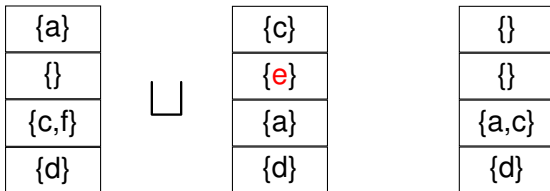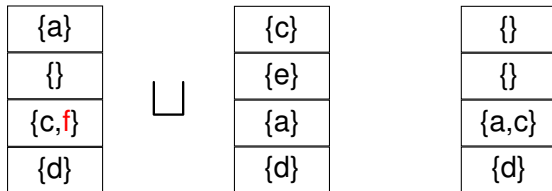| {c} |
|-----|
| {e} |
| {a} |
| {d} |

| {} |
|-----|
| {} |
| {a,c} |
| {d} |

"Intersection + Maximal Age"

# LRU: Must-Analysis: Join

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
- $\gamma(B) \subseteq \gamma(A \sqcup B)$



| {a} |
| {} |
| {c,f} |
| {d} |

$\sqcup$

| {c} |
| {e} |
| {a} |
| {d} |

| {} |
| {} |
| {a,c} |
| {d} |

"Intersection + Maximal Age"

How many memory blocks can be in the must-cache?

entry $[\{\}, \{\}, \{\}, \{\}]$

A $\quad \perp$

$\perp$ B C $\quad \perp$

D $\quad \perp$

exit $\perp$

# Example: Must-Analysis

entry $[\{\}, \{\}, \{\}, \{\}]$

$\perp \sqcup [\{\}, \{\}, \{\}, \{\}] = [\{\}, \{\}, \{\}, \{\}]$

entry $[\{\}, \{\}, \{\}, \{\}]$

$\bot \sqcup [\{\}, \{\}, \{\}, \{\}] = [\{\}, \{\}, \{\}, \{\}]$

A

$[\{A\}, \{\}, \{\}, \{\}]$

B          C          $[\{A\}, \{\}, \{\}, \{\}]$

D  $\bot$

exit  $\bot$

entry  $[\{\}, \{\}, \{\}, \{\}]$

$\perp \sqcup [\{\}, \{\}, \{\}, \{\}] = [\{\}, \{\}, \{\}, \{\}]$

A

$[\{A\}, \{\}, \{\}, \{\}]$

$[\{A\}, \{\}, \{\}, \{\}]$

B          C

D  $[\{B\}, \{A\}, \{\}, \{\}] \sqcup [\{C\}, \{A\}, \{\}, \{\}] =$
$[\{\}, \{A\}, \{\}, \{\}]$

exit  $\perp$

# SAARLAND UNIVERSITY
COMPUTER SCIENCE

entry $[\{\}, \{\}, \{\}, \{\}]$

$[\{D\}, \{\}, \{A\}, \{\}] \sqcup [\{\}, \{\}, \{\}, \{\}] =$
$[\{\}, \{\}, \{\}, \{\}]$

A

$[\{A\}, \{\}, \{\}, \{\}]$

$[\{A\}, \{\}, \{\}, \{\}]$

B

C

D  $[\{B\}, \{A\}, \{\}, \{\}] \sqcup [\{C\}, \{A\}, \{\}, \{\}] =$
$[\{\}, \{A\}, \{\}, \{\}]$

exit  $[\{D\}, \{\}, \{A\}, \{\}]$

No cache hits can be predicted :-(

- Problem:
  - ▶ The first iteration of a loop will always result in cache misses.
  - ▶ Similarly for the first execution of a function.
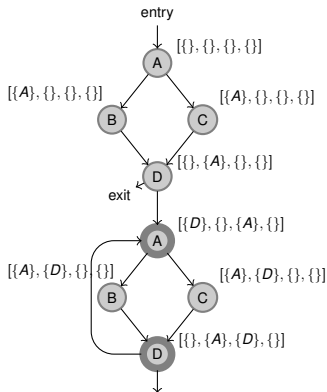- Solution:
  - ▶ Virtually Unroll Loops: Distinguish the first iteration from others
  - ▶ Distinguish function calls by calling context.

Virtually unrolling the loop once:

- Accesses to *A* and *D* are provably hits after the first iteration
- Accesses to *B* and *C* can still not be classified. Within each execution of the loop, they may only miss once.
  $\longrightarrow$ Persistence Analysis

# LRU: May-Analysis: Abstract Domain

- Used to predict *cache misses*.
- Maintains *lower bounds on ages* of memory blocks.
- Lower bound $\geq$ associativity

$\longrightarrow$ memory block definitely *not* cached.

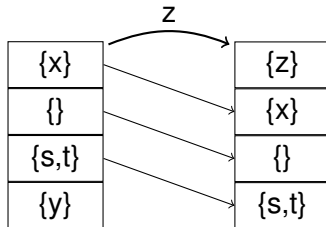## Example

. . . and its interpretation:

Abstract state:

Describes the set of all concrete cache states in which no memory blocks except *x*, *y*, *s*, *t*, and *u* occur,

| | |
|---|---|
| {x,y} | age 0 |
| {} | |
| {s,t} | |
| {u} | age 3 |

- *x* and *y* with an age of at least 0,
- *s* and *t* with an age of at least 2,
- *u* with an age of at least 3.

$\gamma([\{x, y\}, \{\}, \{s, t\}, \{u\}]) =$
$\{[x, y, s, t], [y, x, s, t], [x, y, s, u], \ldots\}$

"Definite Cache Miss":
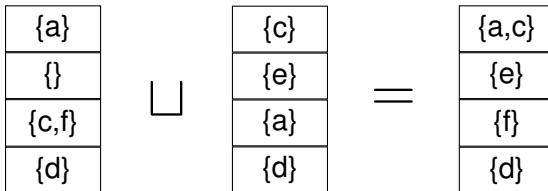
"Potential Cache Hit":

Why does *t* age in the second case?

# LRU: May-Analysis: Join

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
- $\gamma(B) \subseteq \gamma(A \sqcup B)$

| {a} |
|-----|
| {} |
| {c,f} |
| {d} |

$\sqcup$

| {c} |
|-----|
| {e} |
| {a} |
| {d} |

$=$

| {a,c} |
|-------|
| {e} |
| {f} |
| {d} |

"Union + Minimal Age"

# LRU: May-Analysis: Join

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
- $\gamma(B) \subseteq \gamma(A \sqcup B)$

| {a} |
|-----|
| {} |
| {c,f} |
| {d} |

$\sqcup$

| {c} |
|-----|
| {e} |
| {a} |
| {d} |

$=$

| {a,c} |
|-------|
| {e} |
| {f} |
| {d} |

"Union + Minimal Age"

# LRU: May-Analysis: Join

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

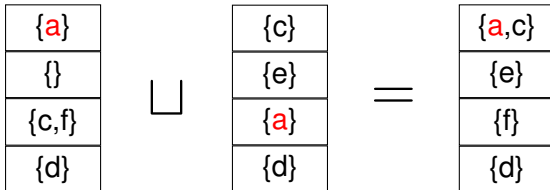- $\gamma(A) \subseteq \gamma(A \sqcup B)$
- $\gamma(B) \subseteq \gamma(A \sqcup B)$

| {a} |
|-----|
| {} |
| {c,f} |
| {d} |

$\sqcup$

| {c} |
|-----|
| {e} |
| {a} |
| {d} |

$=$

| {a,c} |
|-------|
| {e} |
| {f} |
| {d} |

"Union + Minimal Age"

# LRU: May-Analysis: Join

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

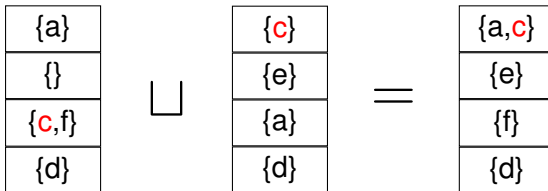- $\gamma(A) \subseteq \gamma(A \sqcup B)$
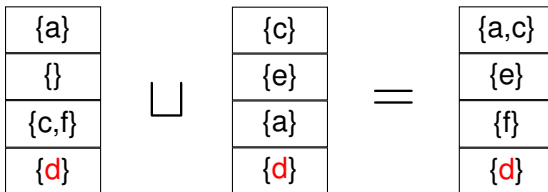- $\gamma(B) \subseteq \gamma(A \sqcup B)$

| {a} |
| --- |
| {} |
| {c,f} |
| {d} |

$\sqcup$

| {c} |
| --- |
| {e} |
| {a} |
| {d} |

$=$

| {a,c} |
| --- |
| {e} |
| {f} |
| {d} |

"Union + Minimal Age"

# LRU: May-Analysis: Join

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
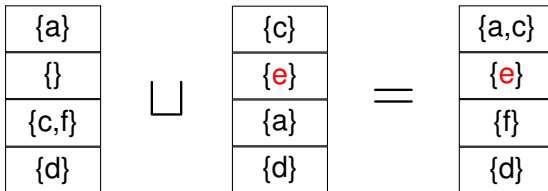- $\gamma(B) \subseteq \gamma(A \sqcup B)$

| {a} |
|-----|
| {} |
| {c,f} |
| {d} |

$\sqcup$

| {c} |
|-----|
| {e} |
| {a} |
| {d} |

$=$

| {a,c} |
|-------|
| {e} |
| {f} |
| {d} |

"Union + Minimal Age"

Need to combine information where control-flow merges.

Join should be conservative (ensures $\gamma$ is monotone):

- $\gamma(A) \subseteq \gamma(A \sqcup B)$
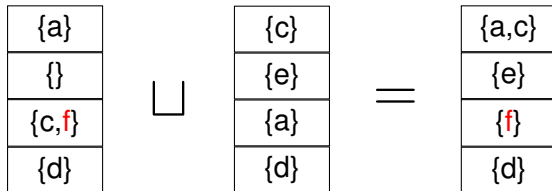- $\gamma(B) \subseteq \gamma(A \sqcup B)$



| {a} |
|-----|
| {} |
| {c,f} |
| {d} |

$\sqcup$

| {c} |
|-----|
| {e} |
| {a} |
| {d} |

=

| {a,c} |
|-------|
| {e} |
| {f} |
| {d} |

"Union + Minimal Age"