

Hash History: A Method for Reconciling Mutual Inconsistency in Optimistic Replication

Brent ByungHoon Kang,
Robert Wilensky and John Kubiatoicz
CS Division, UC Berkeley

Background

- **Optimistic Replication**

- Allow mutable replica to be inconsistent temporarily
 - in a controlled way
 - for high availability and performance
- Tentative update support in OceanStore
- Bayou, USENET, and Peer-to-Peer File System (e.g., Ivy, Pangaea, etc.)

- **Need Mechanism for**

- Figuring out the ordering among updates
- Extracting deltas to be exchanged during reconciliation

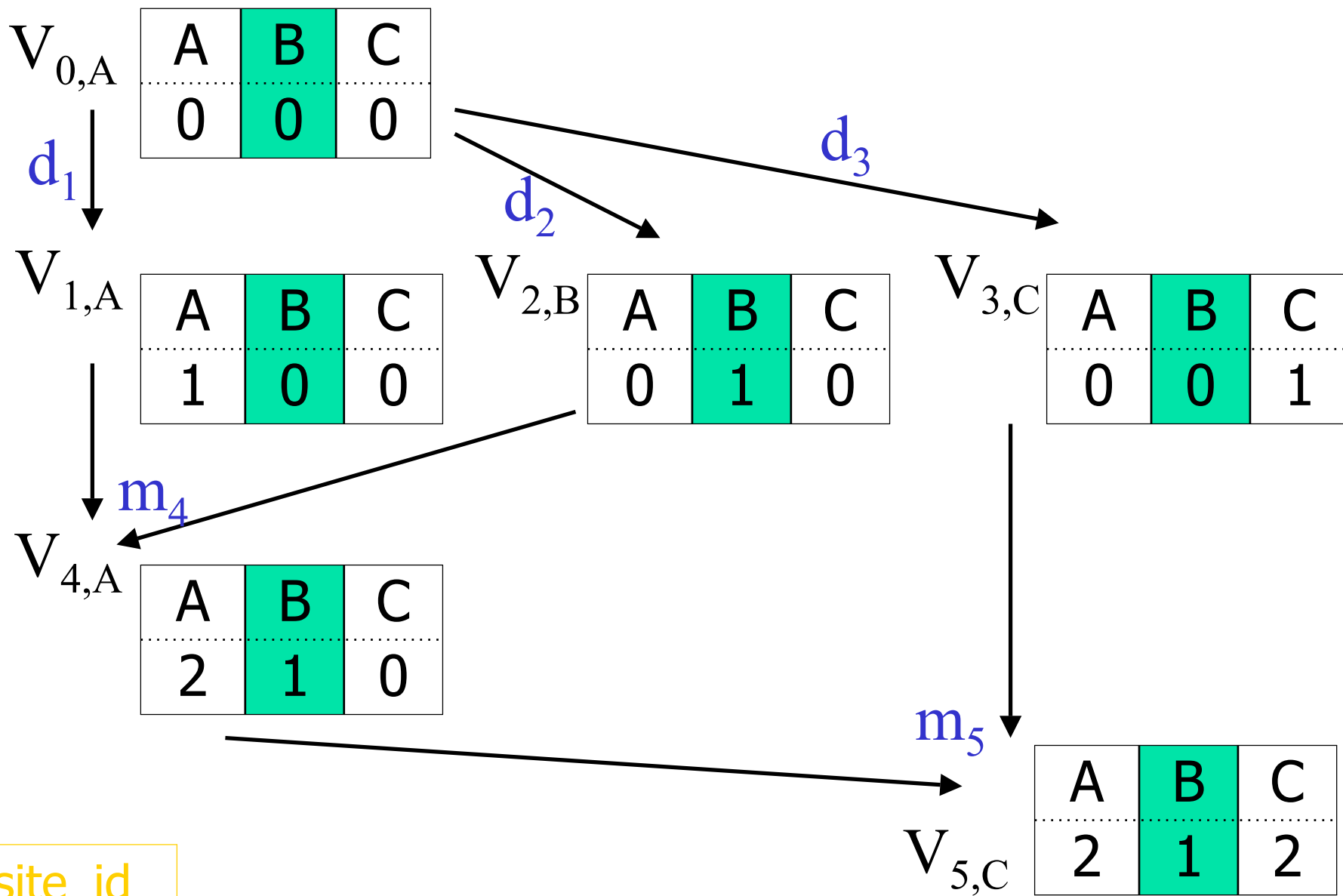
Previous Approaches: Version Vectors

- Widely used in reconciling replicas
 - In most weakly consistent replication systems
 - Bayou, Ficus, Coda, Ivy, Pangaea ... etc.
- Complexity of management grows
 - As new replica site added or deleted
 - Need to assign unique id dynamically for newly added replica sites
- Doesn't scale as number of replica site increases
 - Version vector needs one entry for each replica site
 - Size of vector grows in proportion to number of replica sites

Site A

Site B

Site C



site_id

counter

O1:

A	B	C	D
2	1	0	0

Site A

O1:

A	B	C	D
2	1	0	0

New Site D

O1:

A	B	C	D
0	1	0	0

Site B

O1:

A	B	C	D
2	1	2	0

Site C

Our Proposal: Hash History

- Each site keeps a record of the hash of each version
 - Capture dependency among versions as a directed graph of version hashes (i.e., hash history)
- The sites exchange the hash history in reconciling replicas
- The most recent common ancestral version can be found, if no version dominates
 - Useful hints in a subsequent diffing/merging

$H_{0,A}$

Site A

Site B

Site C

$V_{0,A}$

d_1

$V_{1,A}$

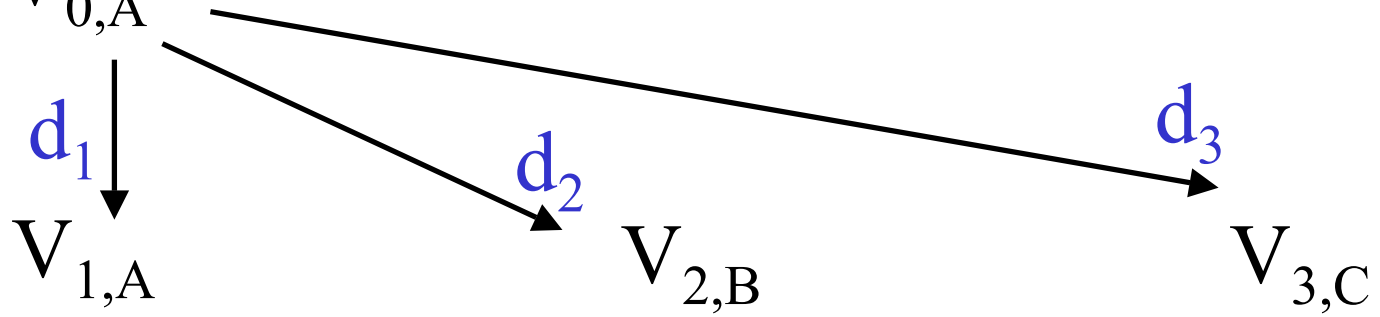
d_2

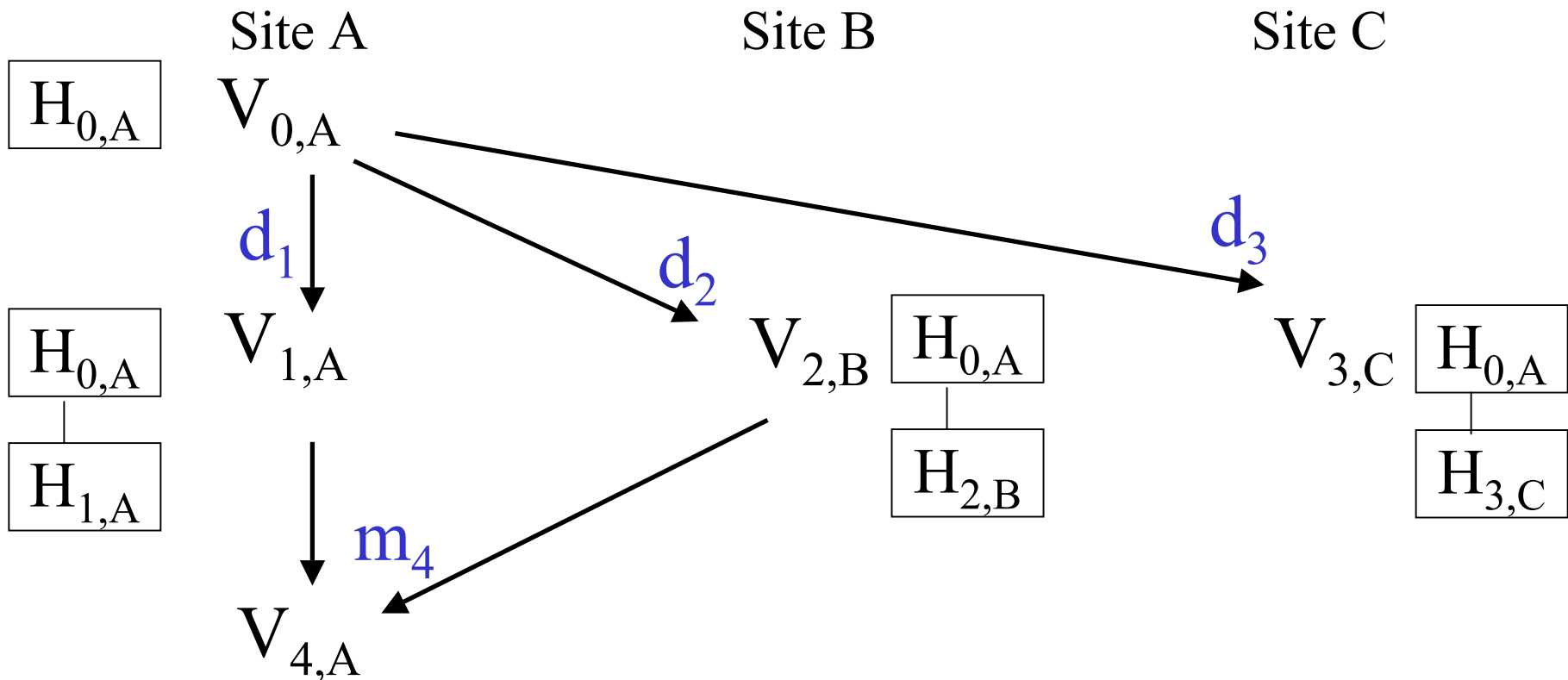
$V_{2,B}$

d_3

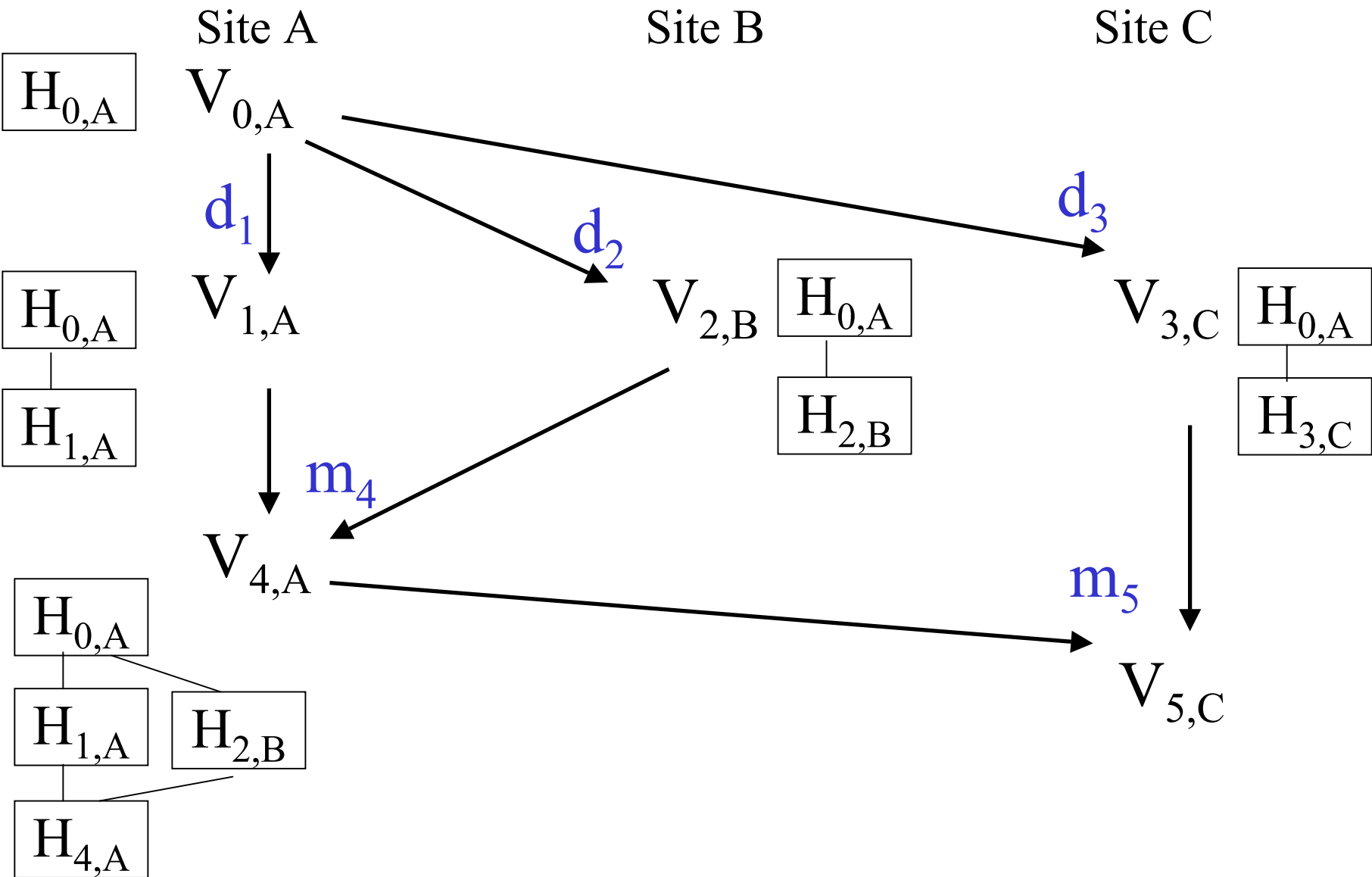
$V_{3,C}$

$$H_{i,site} = \text{hash}(V_{i,site})$$

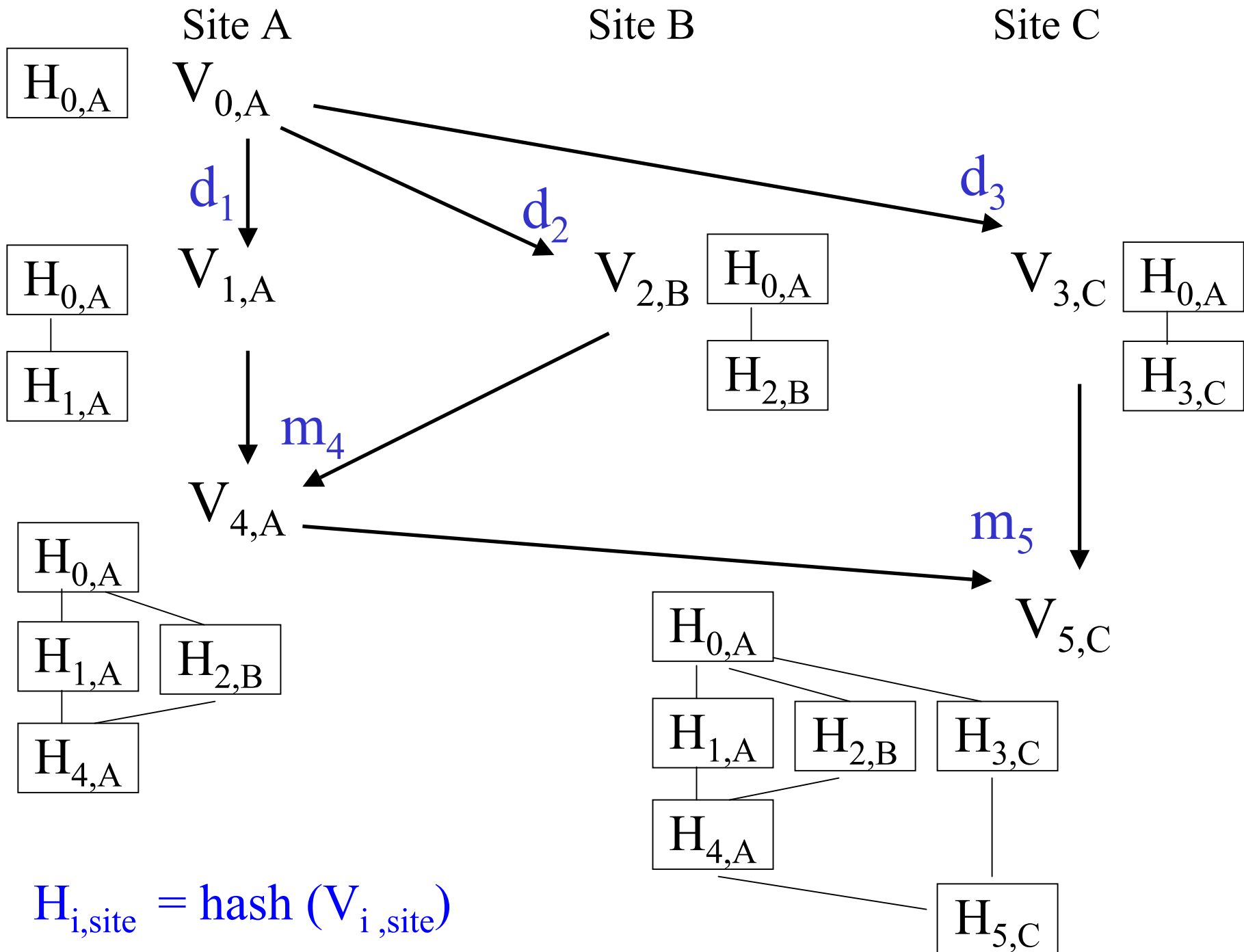




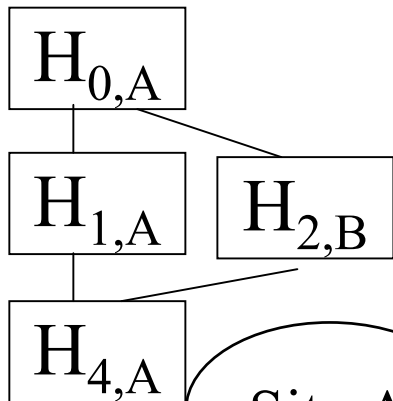
$$H_{i,\text{site}} = \text{hash}(V_{i,\text{site}})$$



$$H_{i,\text{site}} = \text{hash}(V_{i,\text{site}})$$

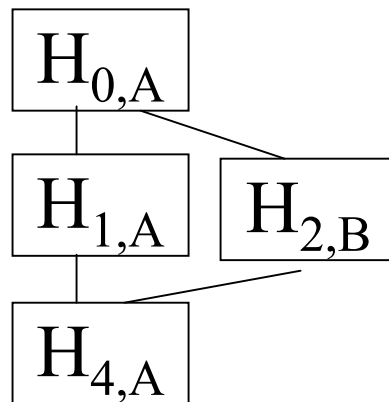


O1:



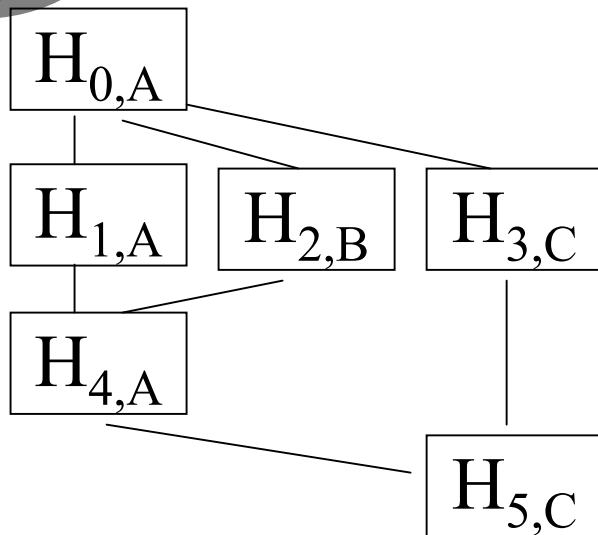
Site A

O1:



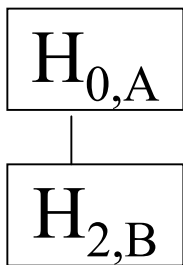
New Site D

O1:

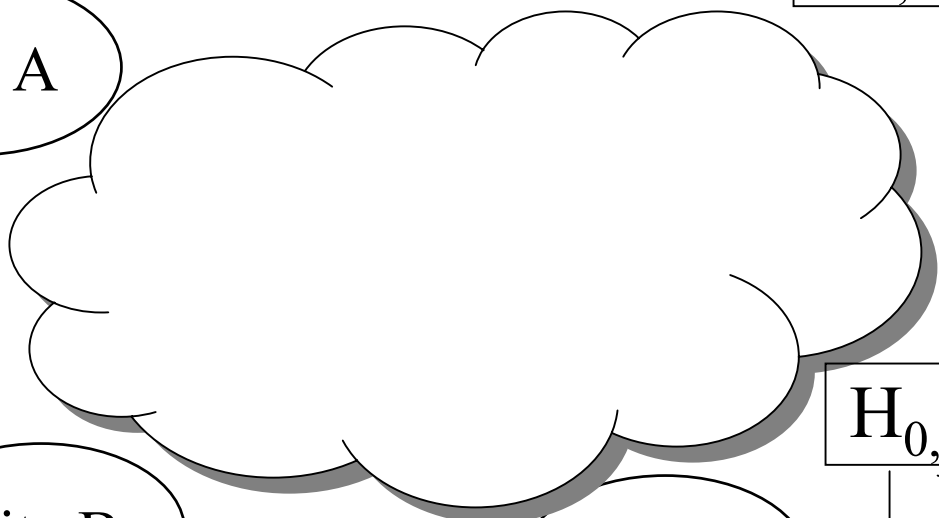


Site C

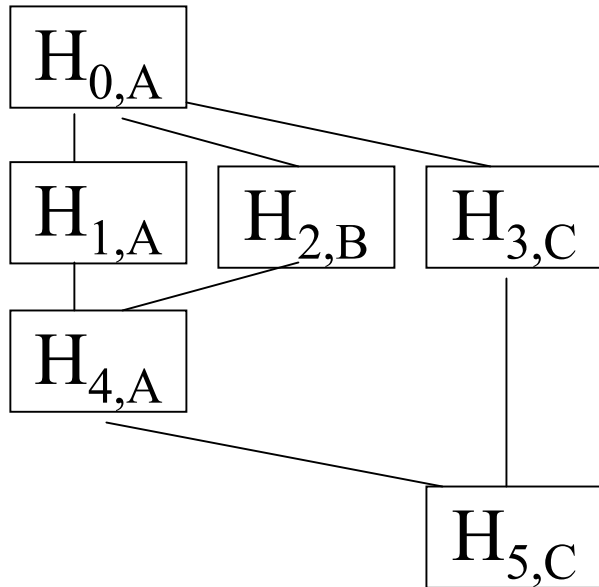
O1:



Site B



Hash History Graph



(a)

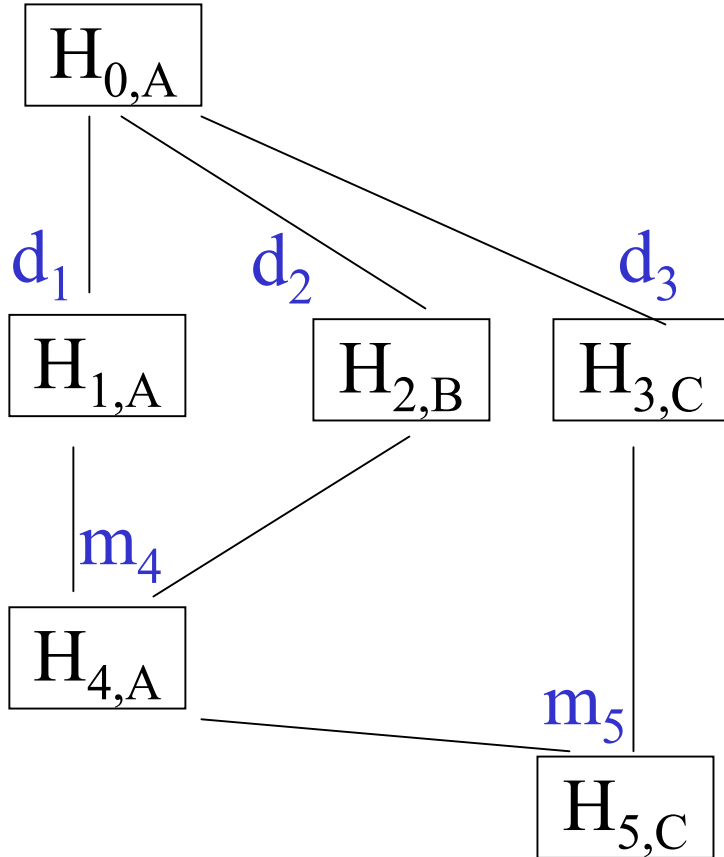
Hash History with Hashtable

Child	Parents
$H_{0,A}$	null
$H_{1,A}$	$H_{0,A}$
$H_{2,B}$	$H_{0,A}$
$H_{3,C}$	$H_{0,A}$
$H_{4,A}$	$H_{1,A} : H_{2,B}$
$H_{5,C}$	$H_{4,A} : H_{3,C}$

Latest : $H_{5,C}$

(b)

Hash History Graph



(a)

Hash History with Hashtable

Child	Parents	delta
$H_{0,A}$	null	null
$H_{1,A}$	$H_{0,A}$	d_1
$H_{2,B}$	$H_{0,A}$	d_2
$H_{3,C}$	$H_{0,A}$	d_3
$H_{4,A}$	$H_{1,A} : H_{2,B}$	m_4
$H_{5,C}$	$H_{4,A} : H_{3,C}$	m_5

Latest : $H_{5,C}$

(b)

HH Properties

- Size of hash history is unbounded
 - Simple Aging
 - Sharable Archived Hash Histories
- Can capture equality case
 - When two different schedule of deltas produce the same output
 - Helps faster convergence

Why Less Conflict in HH than VV

- HH can convey equality information to the descendants while VV cannot
 - E.g., $v1 = \langle A:4, B:5, C:0, D:0, E:0, F:0 \rangle$
 - $v2 = \langle A:5, B:4, C:0, D:0, E:0, F:0 \rangle$
 - C merges then $v3 = \langle A:5, B:5, C:1, D:0, E:0, F:0 \rangle$
 - E merges then $v4 = \langle A:5, B:5, C:0, D:0, E:1, F:0 \rangle$
 - $v3$ and $v4$ could be the same but VV shows conflict !
- If $v3$ and $v4$ are considered equal, then
 - all descendants of $v4$ will dominate $v3$.
- If $v3$ and $v4$ are considered as in conflict,
 - all descendants of $v4$, will be in conflict with $v3$

Experiment Goal

- Comparison with version vector result:
 - HH converges faster with a lower conflict rate than a version vector scheme
 - To what extent is this true in practice?
- Aging Policy:
 - the aging period for pruning hash history
 - vs. HH size
 - vs. the false conflict rate due to aging
 - when the pruned part of the hash history is required for determining the version dominance

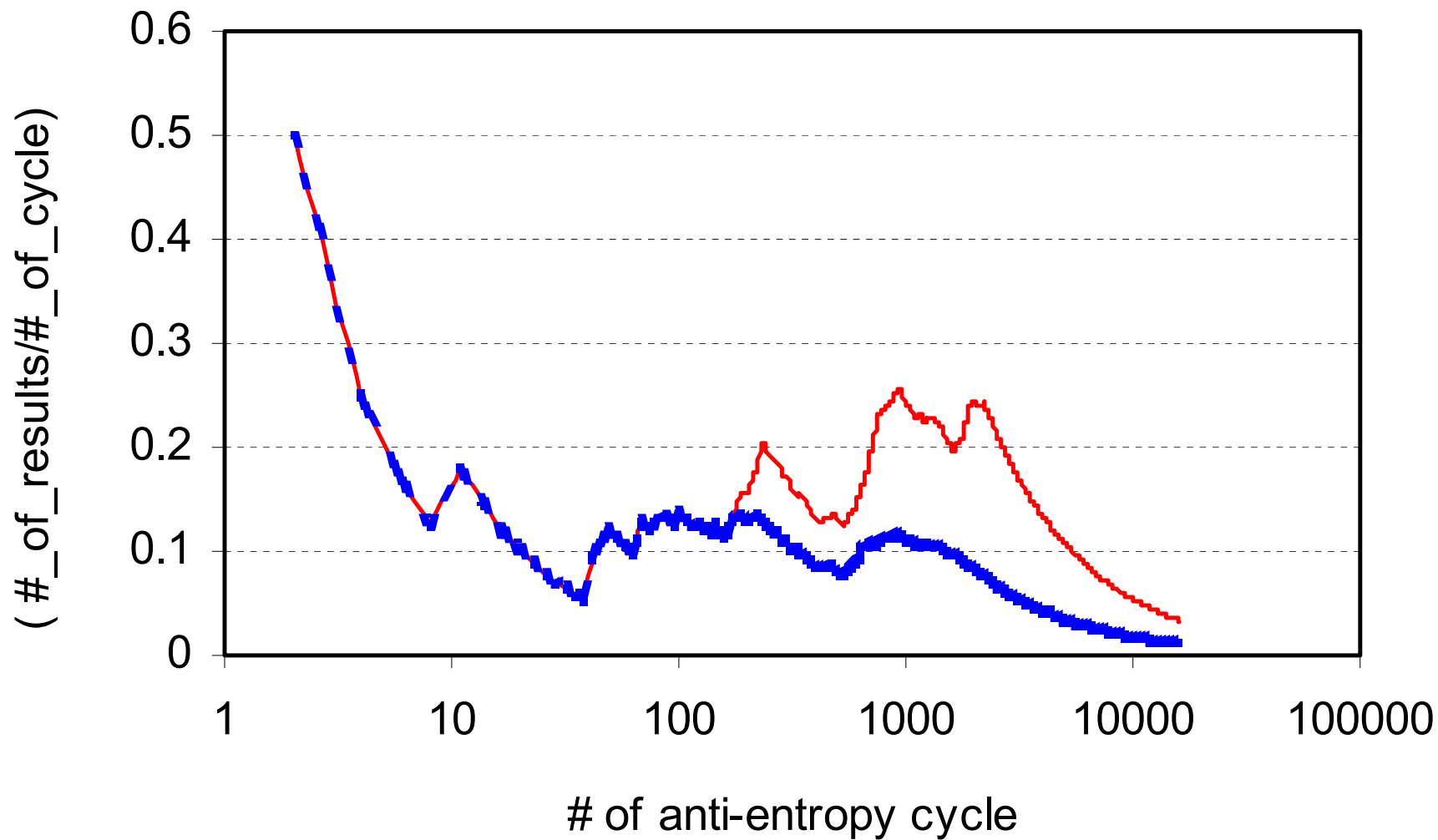
Simulation Setup

- Event-driven simulator
 - Events are collected from CVS logs
 - Each user represent a replica site
 - Reads the event `<time, user, filename>`
 - After each event, the simulator
 - repeats the anti-entropy for 50% (or 25%) of the total number of sites.
 - E.g., if there are 20 sites so far, the anti-entropy is repeated for 10 times with 50% parameter after each event.

CVS Trace Data (from sourceforge.net)

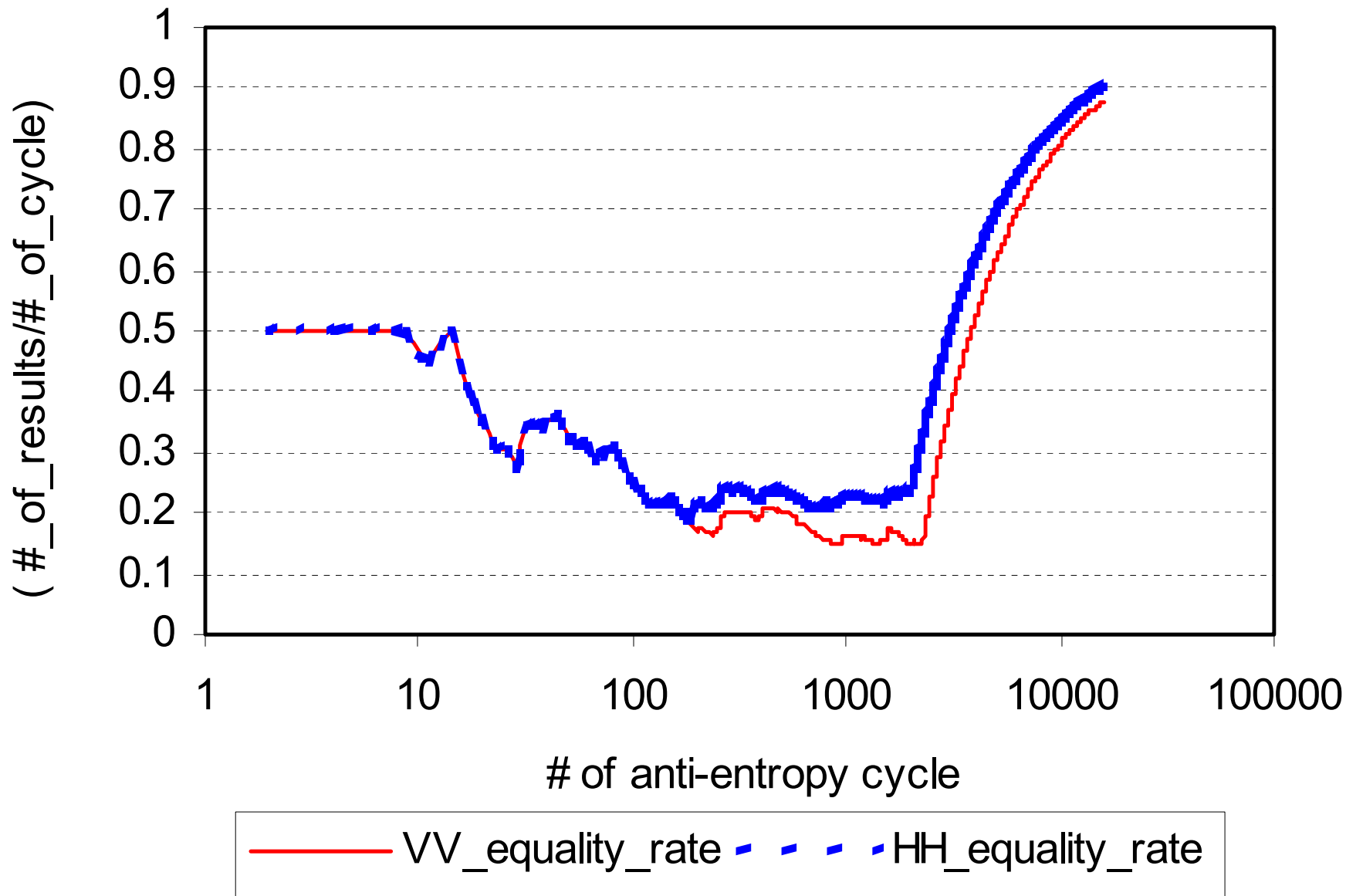
	Dri	Freenet	Pcgen
# of events	10137	2281	404
# of users	21	64	39
Duration	4/27/1994 - 5/3/2002	12.28.1999 -4/25/2002	1/17/2002 - 4/12/2002
inter-commit time AVG	101.3 min	237.8 min	225.4 min
MEDIAN	0.016 min	34.6 min	2.16 min

Conflict rate of VV and HH

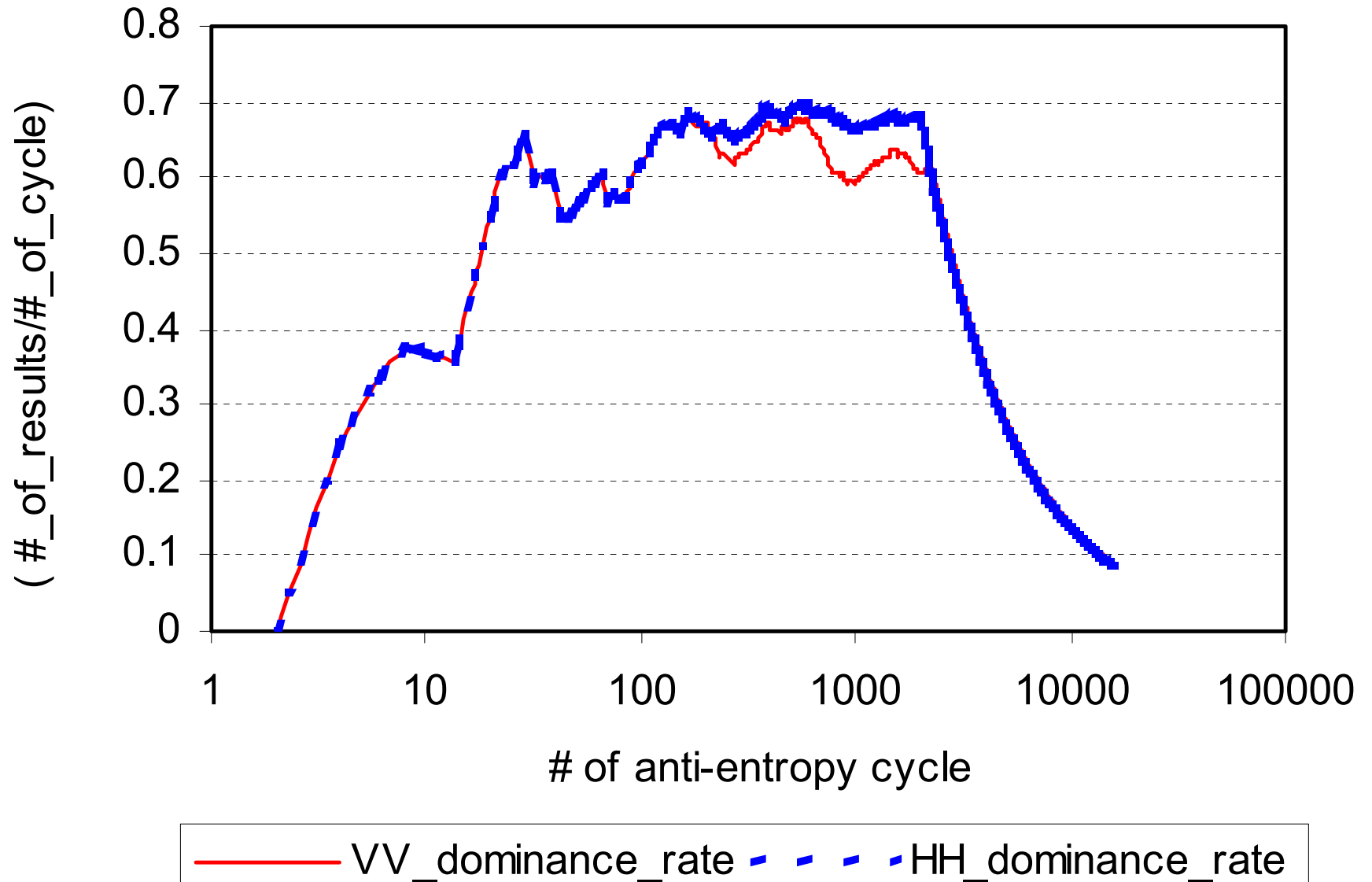


— VV_conflict_rate - - - HH_conflict_rate

Equality rate of VV and HH



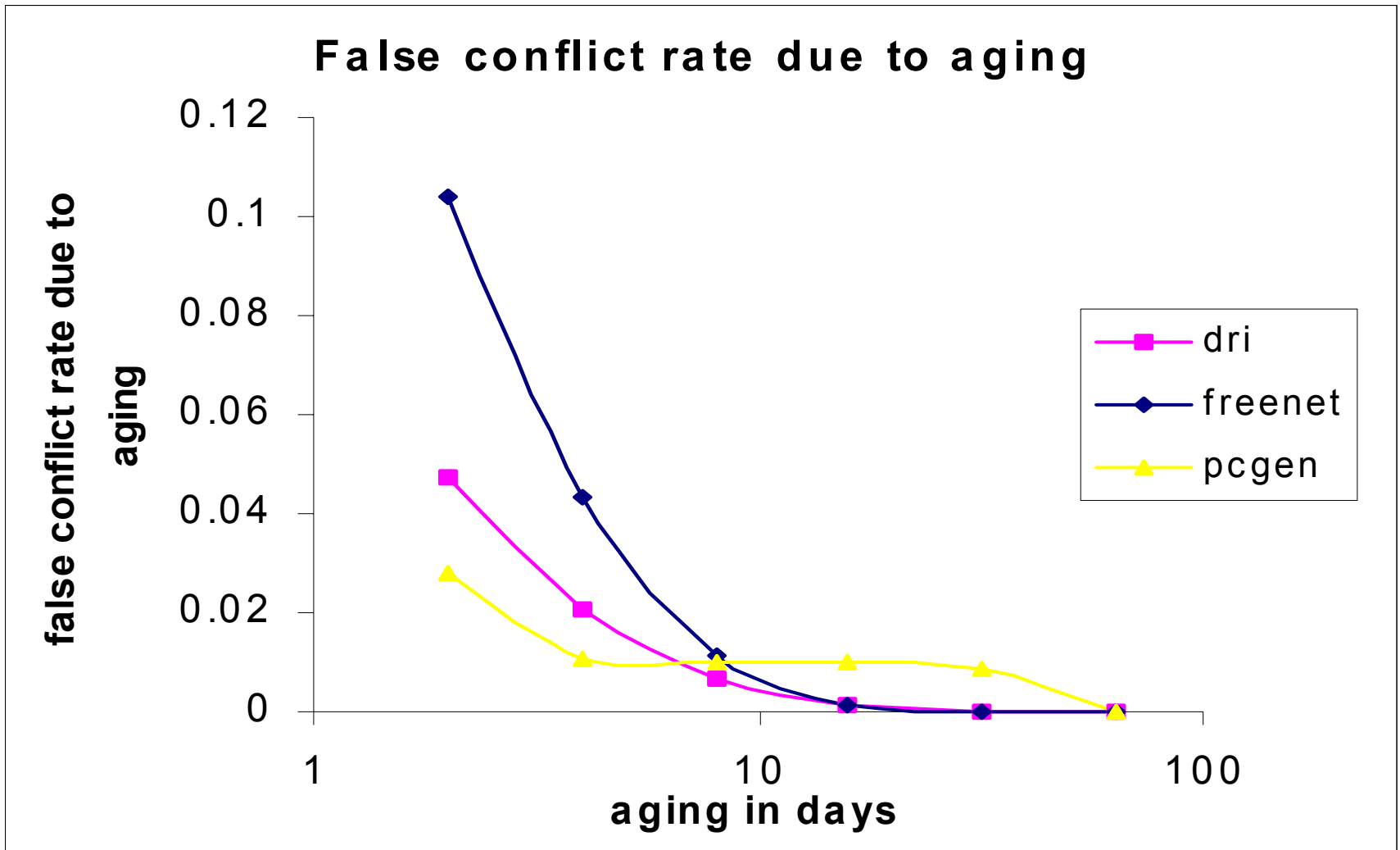
Dominance rate of VV and HH



Aging Period vs. HH Size

Aging period (days)	HH size (# of entries) – dri	pcgen	freenet	Average
32	146.3	159.1	61.5	122.3
64	413.9	443.9	147.5	335.1
128	551.5	591.7	612.8	585.3

Aging Period vs. False Conflict



Conclusion

- Simple to maintain
 - No complexity in site addition/deletion
 - No need to assign unique id dynamically for newly added replica sites
- Scalable to thousands of sites
 - HH grows in proportion to number of update instances not number of sites
- Faster Convergence
 - HH can capture and propagate equality information
- HH growth can be controlled effectively by
 - using aging policy or sharing archived hash history

Future Work

- Security aspect of HH
 - Self-verifiable
 - Can detect mal-functioning site
- More information
 - Hash History Approach for Reconciling Mutual Inconsistency in Optimistic Replication, *B. Kang, R. Wilensky and J. Kubiatoicz*, The 23rd International Conference on Distributed Computing Systems (ICDCS), 2003, Providence, Rhode Island USA
 - <http://www.cs.berkeley.edu/~hoon/hashhistory>

Site A

Site B

Site C

$$V_{0,A}$$

A	B	C
0	0	0
σ_A^0	σ_B^0	σ_C^0

$$V_{1,A}$$

A	B	C
1	0	0
σ_A^1	σ_B^0	σ_C^0

$$V_{4,A}$$

A	B	C
2	1	0
σ_A^2	σ_B^1	σ_C^0

$$V_{2,B}$$

A	B	C
0	1	0
σ_A^0	σ_B^1	σ_C^0

$$V_{3,C}$$

A	B	C
0	0	1
σ_A^0	σ_B^0	σ_C^1

$$V_{5,C}$$

A	B	C
2	1	2
σ_A^2	σ_A^1	σ_A^2

site
counter
signature

$\sigma_{\text{site}}^{\text{value}} = \text{sign}_{\text{site}}(\text{value of site's local counter})$

Site A

Site B

Site C

$$V_{0,A}$$

nil
H_0
α_A^0

$$V_{1,A}$$

nil	H_0
H_0	H_1
α_A^0	α_A^1

$$V_{2,B}$$

nil	H_0
H_0	H_2
α_A^0	α_B^2

$$V_{3,C}$$

nil	H_0
H_0	H_3
α_A^0	α_C^3

$$V_{4,A}$$

nil	H_0	H_0	$H_1 : H_2$
H_0	H_1	H_2	H_4
α_A^0	α_A^1	α_B^2	α_A^4

$$V_{5,C}$$

nil	H_0	H_0	$H_1 : H_2$	H_0	$H_3 : H_4$
H_0	H_1	H_2	H_4	H_3	H_5
α_A^0	α_A^1	α_B^2	α_A^4	α_C^3	α_C^5

parents
child
authentic
icator

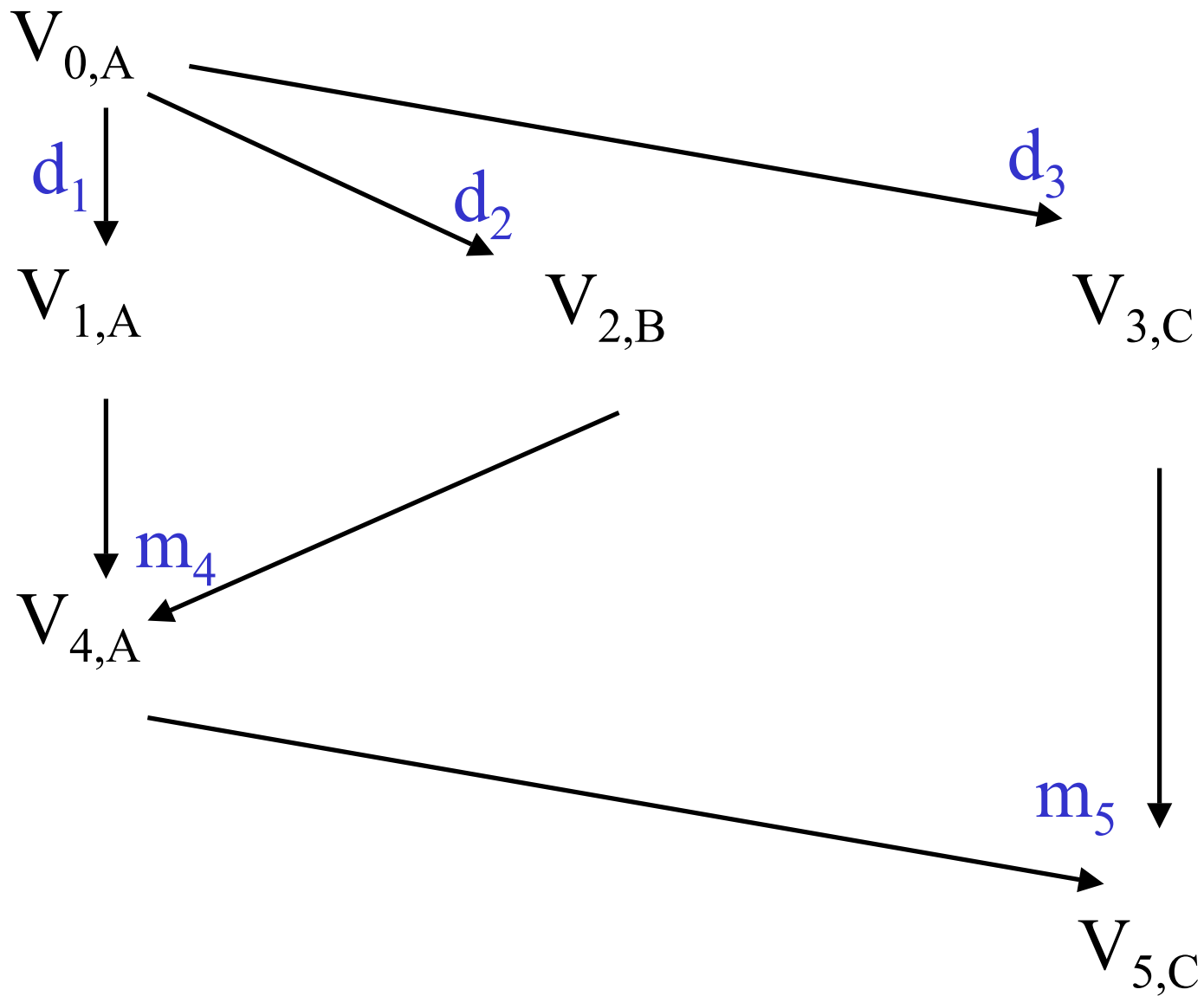
$$H_i = \text{hash}(V_{i,\text{site}})$$

$$\alpha_{\text{site}}^k = \text{hash}(\alpha^k\text{'s parents} \parallel H_k)$$

Site A

Site B

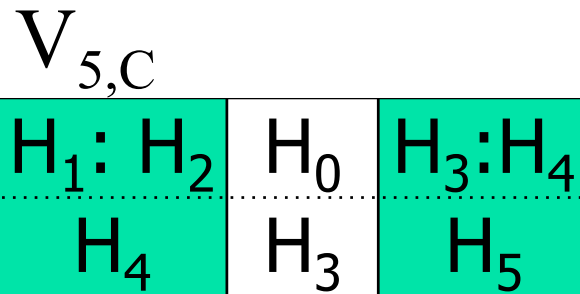
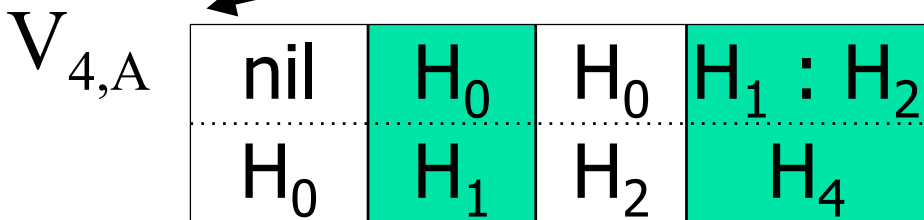
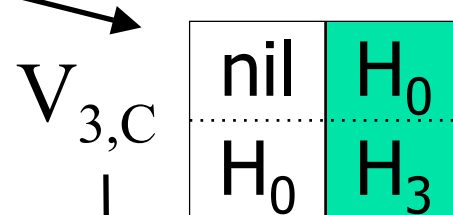
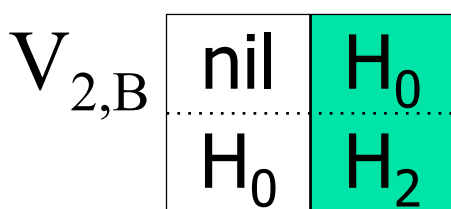
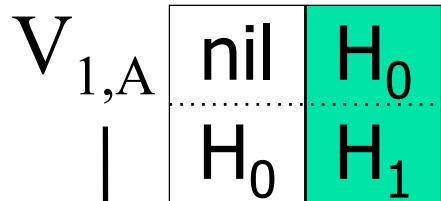
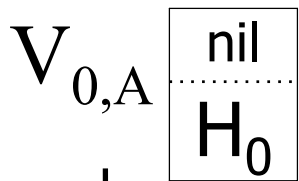
Site C



Site A

Site B

Site C



parents
child

$H_i = \text{hash}(V_{i,\text{site}})$