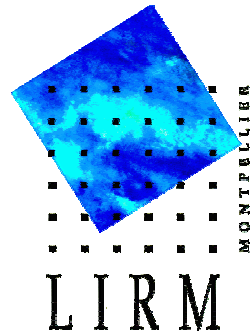

Asynchronous Reconciliation

based on Operational Transformation

for P2P Collaborative Environments

Michelle CART

Jean FERRIÉ



Université Montpellier II
Montpellier (France)

The Context

Reconciling multiple distributed copies of an object

Synchronous methods:

- copies reconciled in real-time
- need broadcast protocol
- based on Operational Transformation (CSCW)
- need ordering mechanisms (state vectors, timestamps, sequencer, ..)

➔ unsuitable for P2P environments

Asynchronous methods:

- pair-wise copy reconciliation
- need central component:
 - master copy (Copy/ Modify/ Merge paradigm)
 - particular site [Bayou, IceCube], sequencer [SO6]

➔ unsuitable for P2P environments

Our Proposal

Multiple divergent copies of an object:



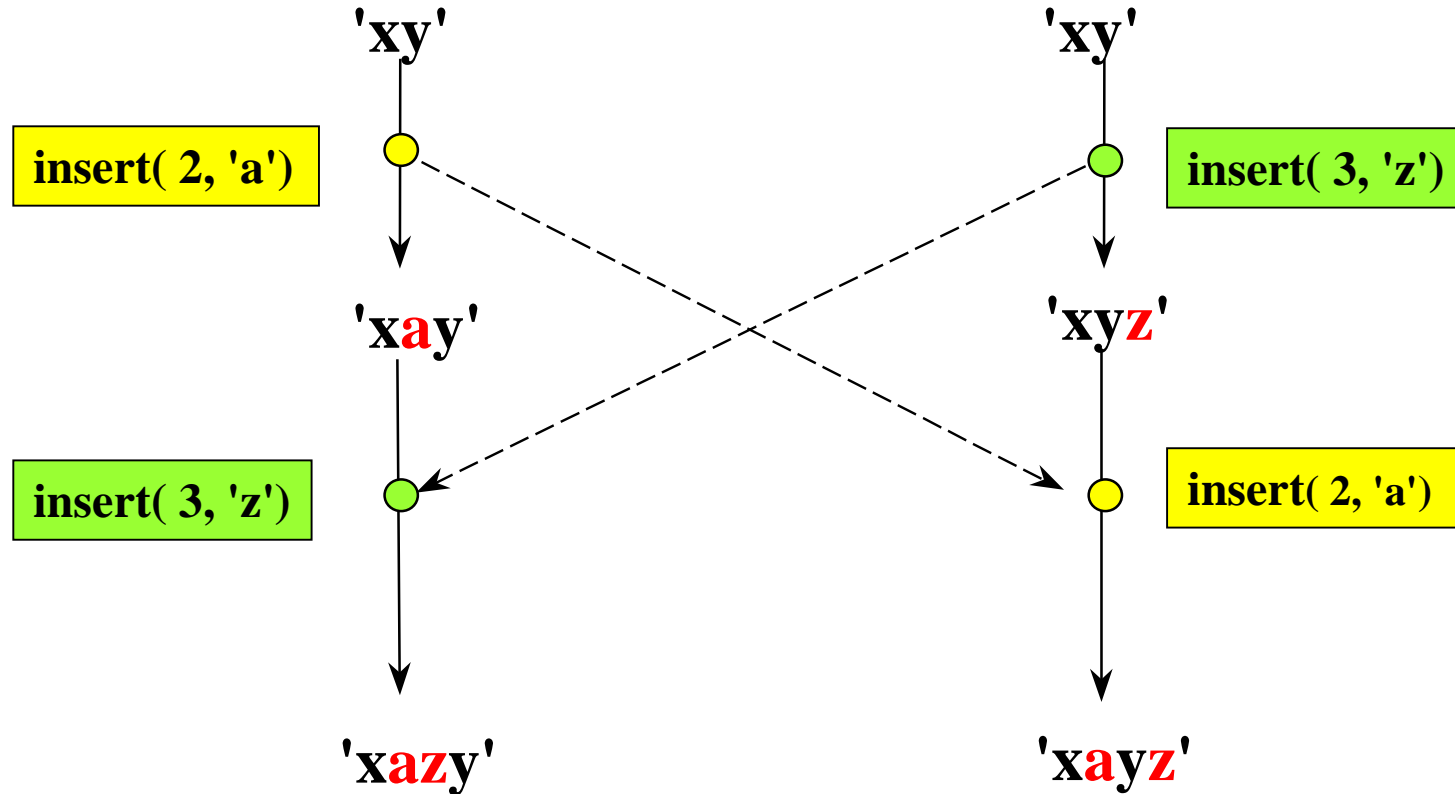
Reconciliation

- ensuring "eventual consistency" of the copies
- without rejecting any operation
- asynchronous (at any time)
- pair-wise (regardless of the pair, no privileged copy)
- group of an undetermined size
(sites can dynamically join / leave the group)
- scalable to a P2P network
- based on Operational Transformation

Operational Transformation: recall



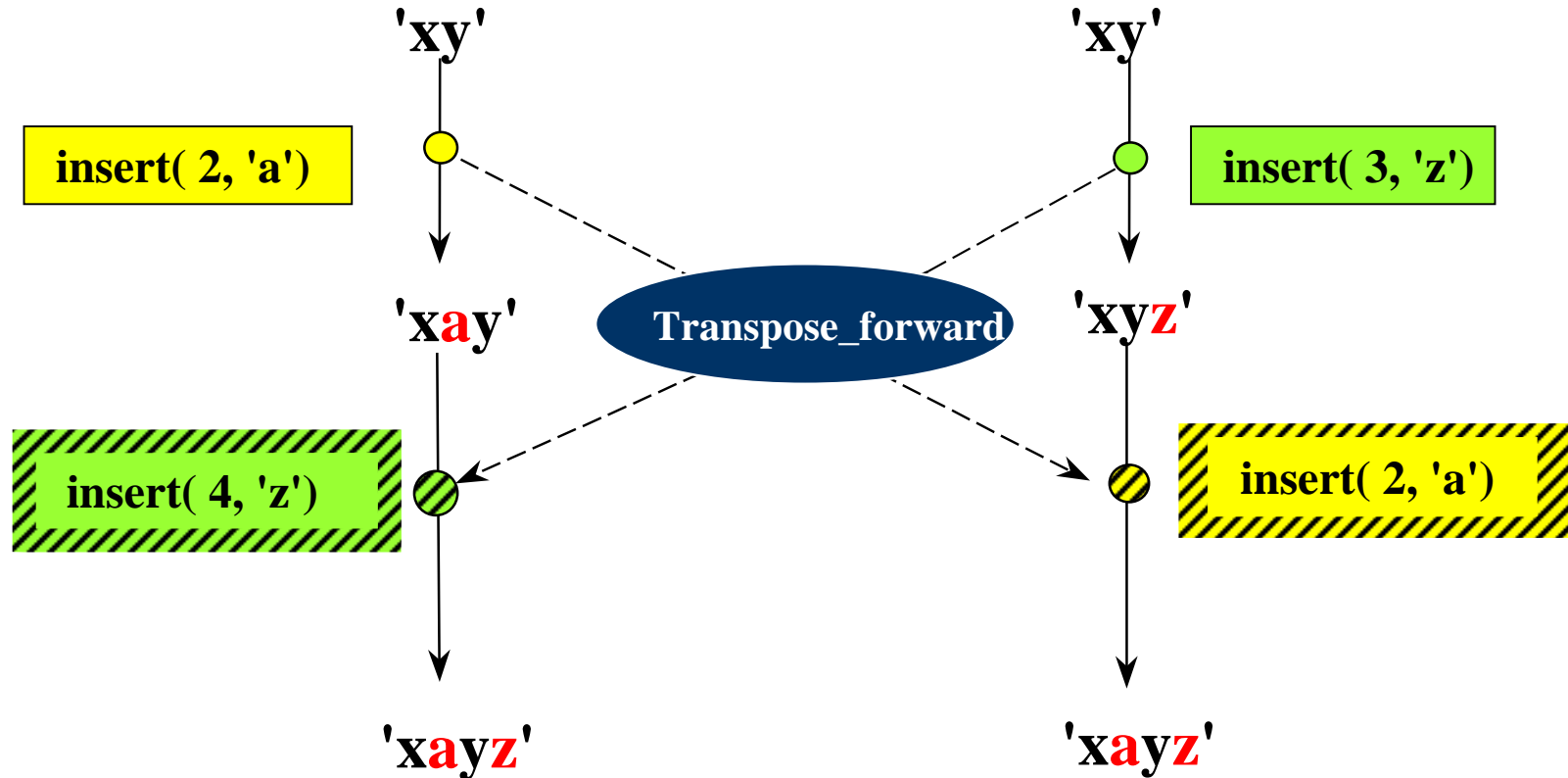
Used to achieve copy convergence when concurrent and non-commutative operations are executed on copies in different orders



Operational Transformation: recall



Used to achieve copy convergence when concurrent and non-commutative operations are executed on copies in different orders



Transformation functions specified (before execution)
Conflicting situations automatically solved (during execution)

Operational Transformation: example

`Transpose_forward (insert(p1, c1) , insert(p2, c2)) =`

`case p1 ? p2 of`

`p1 < p2 : return insert(p2 +1, c2);`

`p1 > p2 : return insert(p2, c2);`

`p1 = p2 : if c1 = c2 then return id`

`else`

`if pr(c2) > pr(c1)`

`then return insert(p2, c2)`

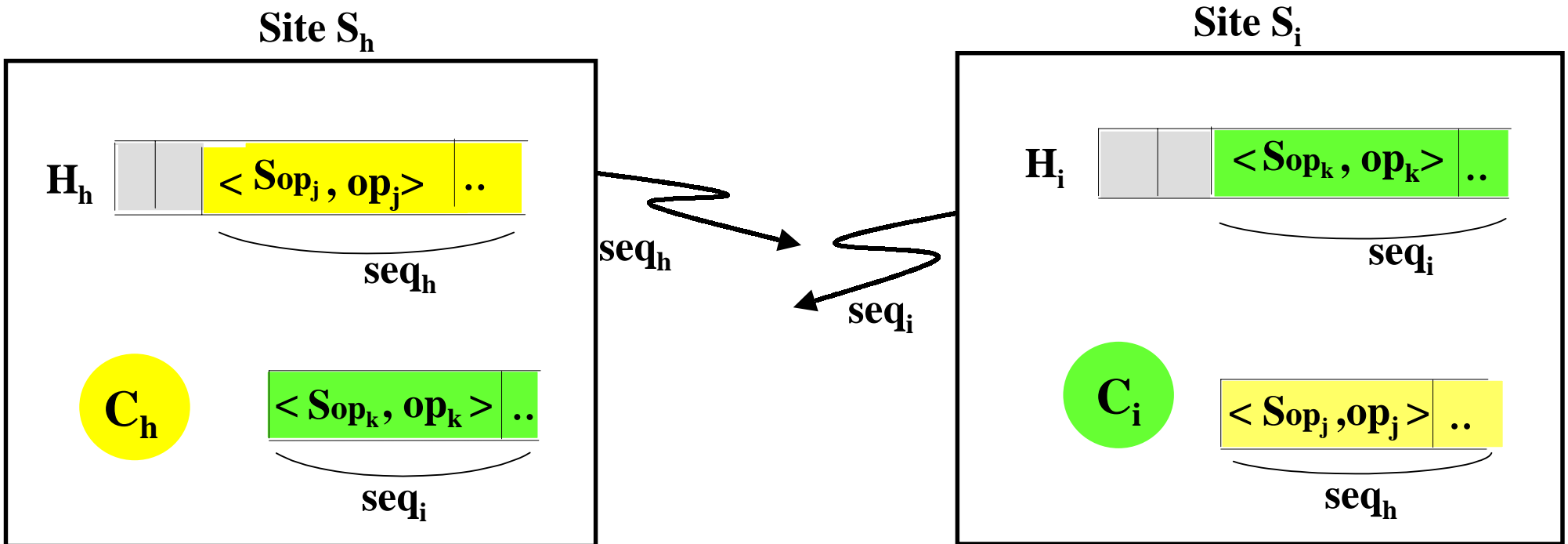
`else return insert(p2+1, c2);`

`endif;`

`endif;`

`endcase`

The Model (1)



History H : succession of operations $\langle S_{op}, op \rangle$ that leads from the initial state to the current state of the copy C .

\uparrow \leftarrow operation executed on the copy C associated to H
 generator site of op

Reconciliation \rightarrow transmission of sequences seq of operations: $seq_h \rightarrow S_i$
 $seq_i \rightarrow S_h$

The Model (2)

Relationships between operations $\langle S_{op_j}, op_j \rangle$ and $\langle S_{op_k}, op_k \rangle$

→ **Causal precedence** "happen before"

$$op_j \xrightarrow{c} op_k \left\{ \begin{array}{l} S_{op_j} = S_{op_k} \text{ and } op_j \text{ generated before } op_k \\ S_{op_j} \neq S_{op_k} \text{ and } op_k \text{ generated after the execution of } op_j \text{ on } S_{op_k} \\ \exists op_m : (op_j \xrightarrow{c} op_m) \text{ and } (op_m \xrightarrow{c} op_k) \end{array} \right.$$

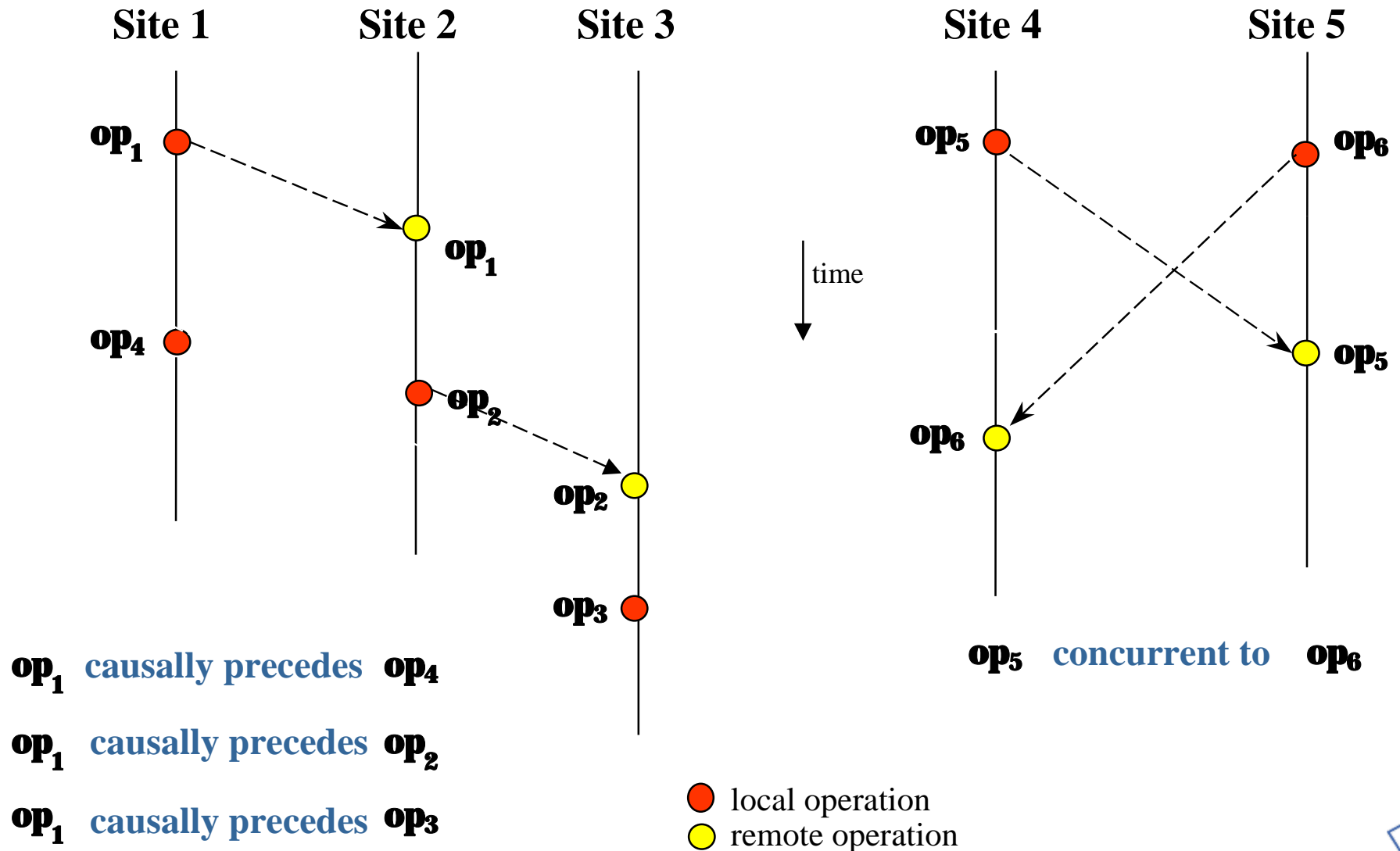
→ **Concurrency**

$$op_j // op_k \quad \text{not } (op_j \xrightarrow{c} op_k) \text{ and not } (op_k \xrightarrow{c} op_j)$$

→ **Precedence**

$$op_j \xrightarrow{H} op_k \quad op_j \text{ appears before } op_k \text{ in history } H$$

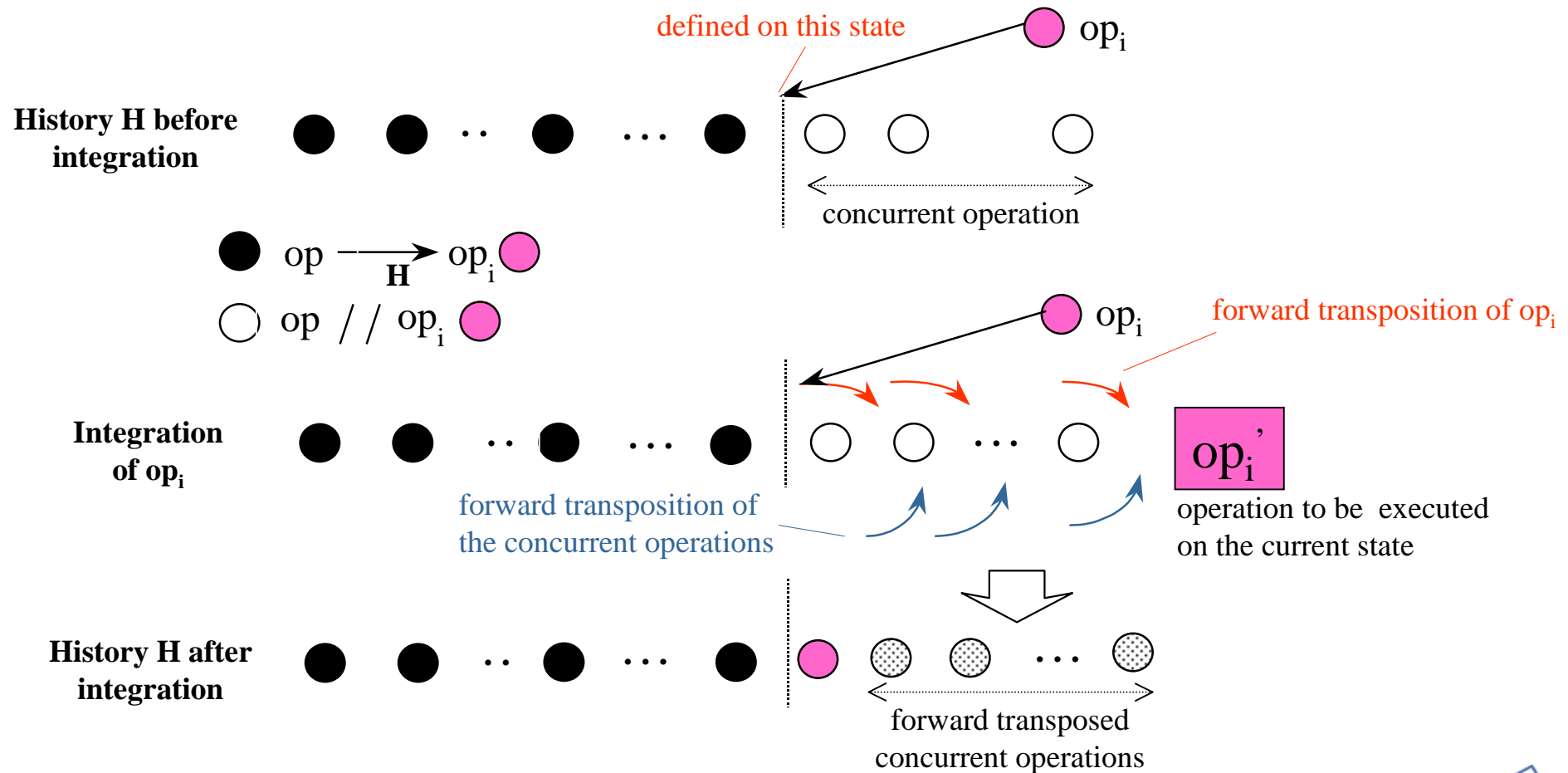
The model (2): example



Integration of an operation

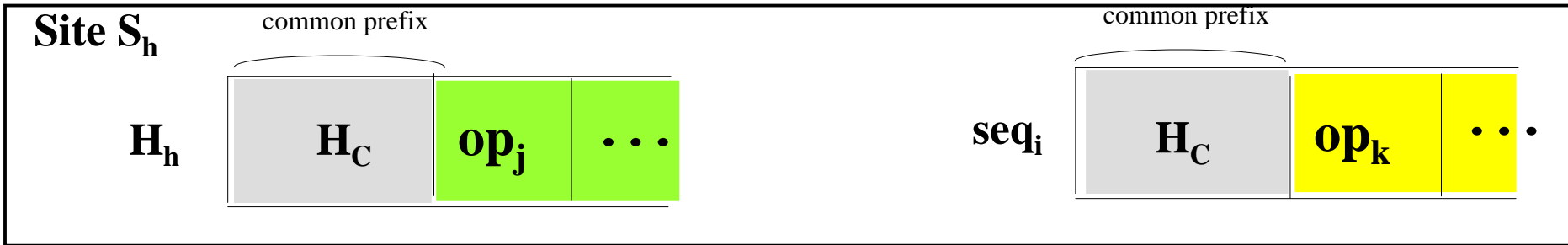


Integration of the operation op_i into the history H of a site

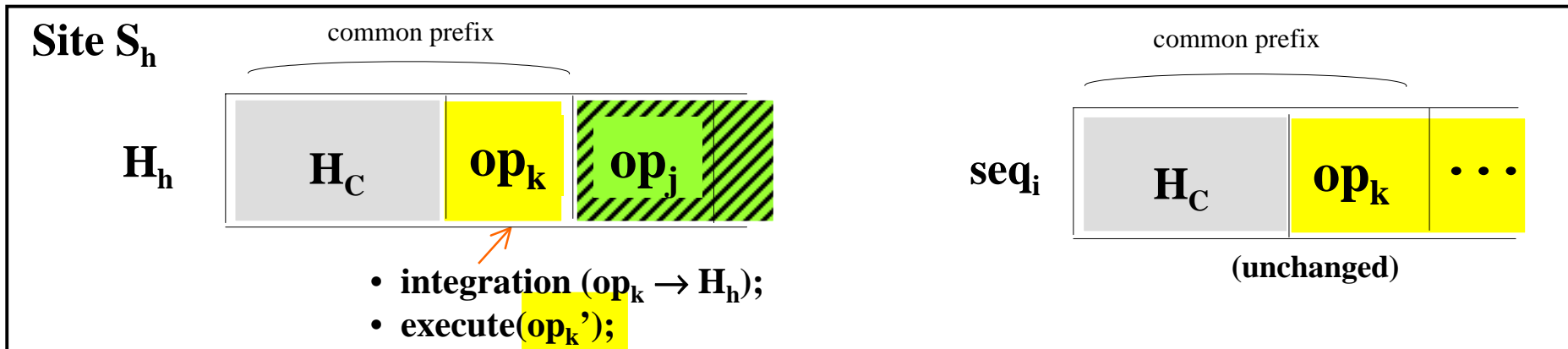
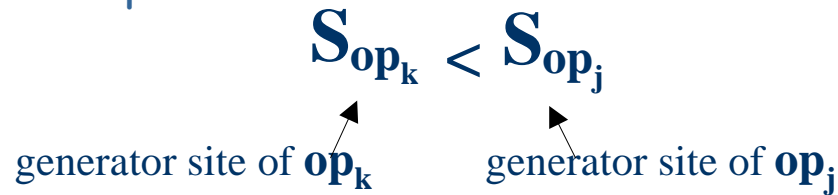


Basic Principle of MOT2

("Merge based on Operational Transformation")

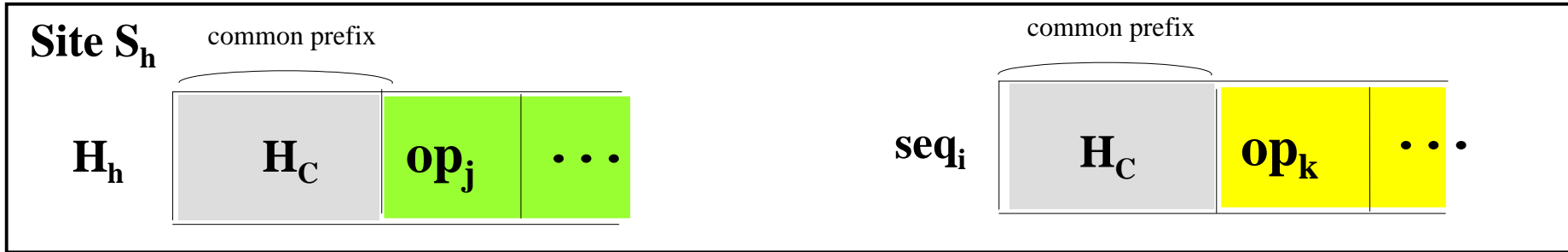


The order between 2 operations to be merged depends on the **order defined among the generator sites** of these operations.



Basic Principle of MOT2

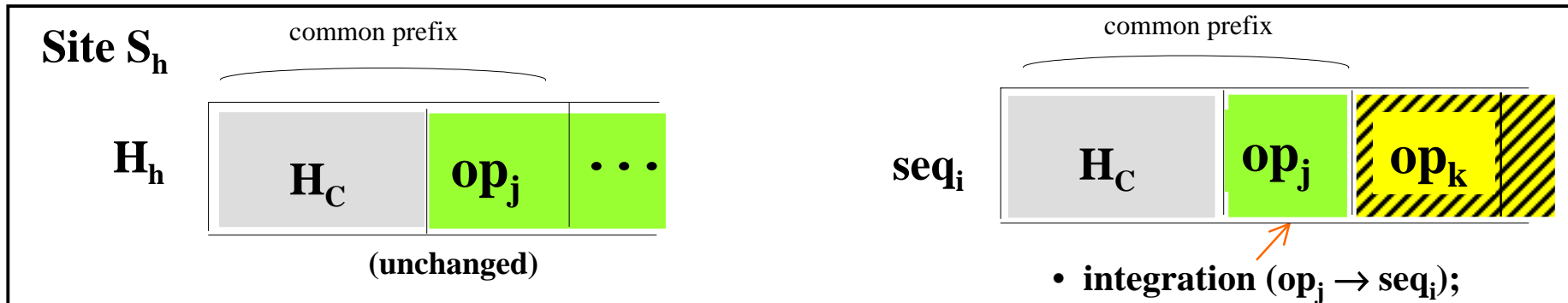
("Merge based on Operational Transformation")



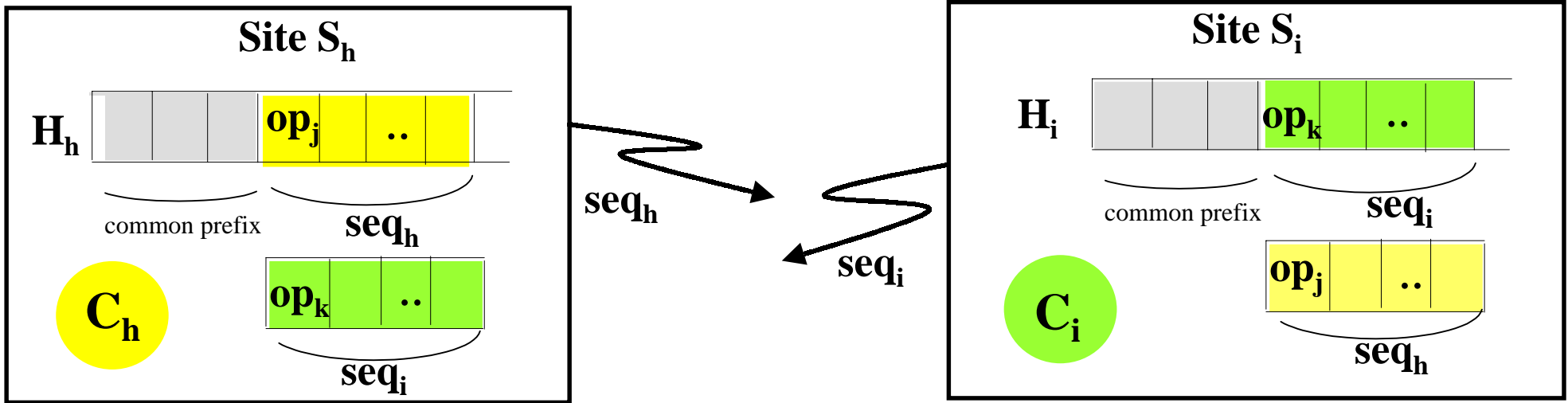
The order between 2 operations to be merged depends on the **order defined among the generator sites** of these operations.

$$S_{op_j} < S_{op_k}$$

\swarrow generator site of op_j \swarrow generator site of op_k



The MOT2 Merge Algorithm



while loop

..

case $S_j ? S_k$ of

$S_k < S_j$: integration ($op_k \rightarrow H_h$) ; execute(op_k);

$S_j < S_k$: integration ($op_j \rightarrow seq_i$);

$S_k = S_j$: ; -- $op_j = op_k$

end loop ;

while loop

..

case $S_j ? S_k$ of

$S_k < S_j$: integration ($op_k \rightarrow seq_h$);

$S_j < S_k$: integration ($op_j \rightarrow H_i$) ; execute(op_j);





$S_k = S_j$: ; -- $op_j = op_k$


end loop ;


Algorithm symmetrical and generic

The MOT2 Algorithm

```

procedure MOT2 ( $H_k, H_j$ );
{Look for the prefix  $H_C$  common to  $H_k$  and  $H_j$ : détermine the index  $i_s$  of the last operation of  $H_C$ }
   $i := i_s + 1$ ;
while ( $i \leq \text{size\_}H_k$ ) and ( $i \leq \text{size\_}H_j$ ) loop
   $\langle \text{Idop}_k, \text{Sop}_k, \text{op}_k \rangle := H_k[i]$ ; 
   $\langle \text{Idop}_j, \text{Sop}_j, \text{op}_j \rangle := H_j[i]$ ; 
  case  $\text{Sop}_k ? \text{Sop}_j$  of
     $\text{Sop}_k < \text{Sop}_j$ : Integration ( $H_j, i, \langle \text{Idop}_k, \text{Sop}_k, \text{op}_k \rangle$ );  $\{ \text{op}_k \text{ integrated into } H_j \}$  
     $\text{Sop}_j < \text{Sop}_k$ : Integration ( $H_k, i, \langle \text{Idop}_j, \text{Sop}_j, \text{op}_j \rangle$ );  $\{ \text{op}_j \text{ integrated into } H_k \}$  
     $\text{Sop}_k = \text{Sop}_j$ : ;  $\{ \text{Operation already present in } H_k \text{ and } H_j \}$ 
  endcase;
   $i := i + 1$ ;
endloop;

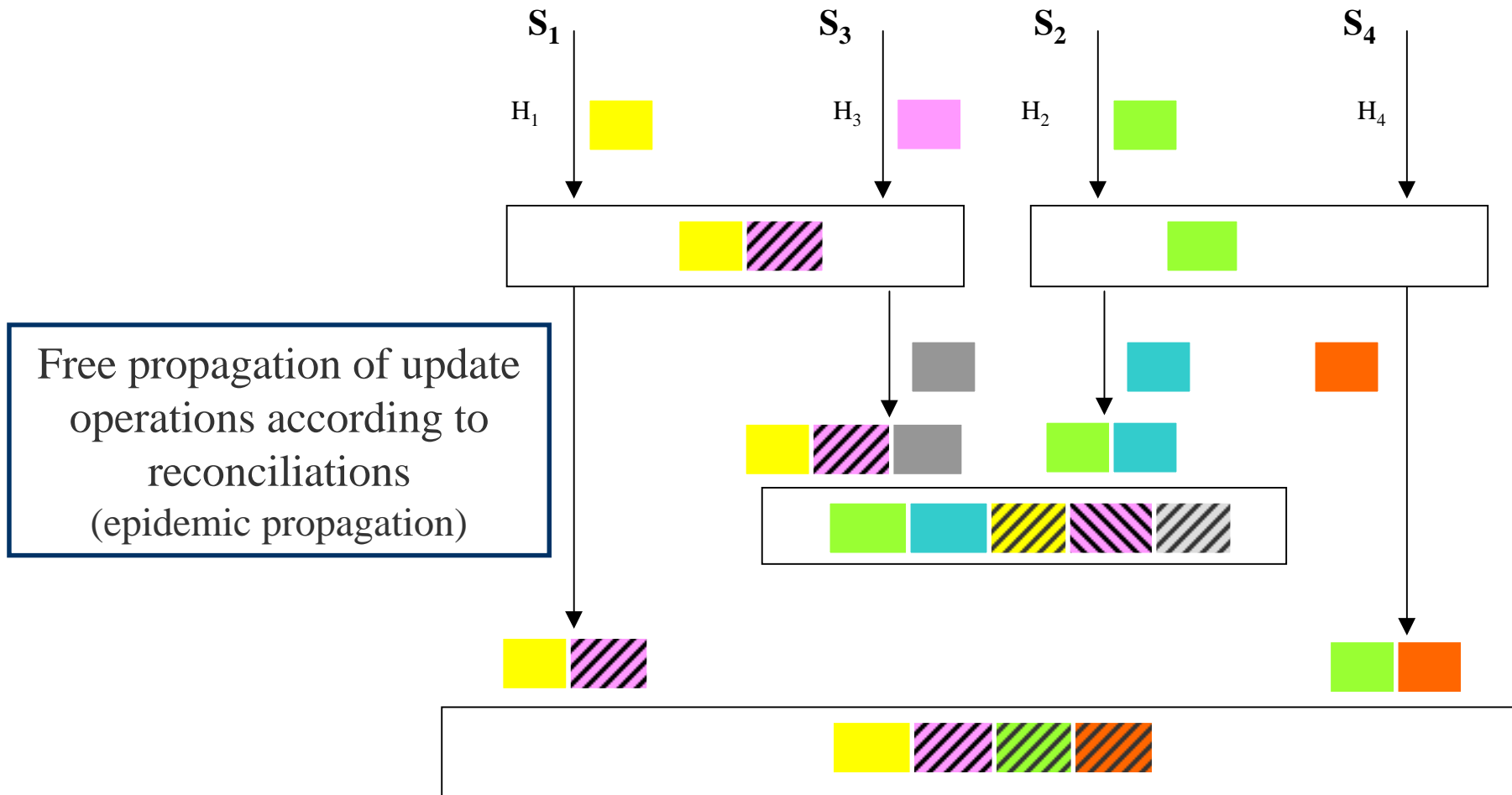
  {The end of  $H_k$  or  $H_j$  has been reached}
  while  $i \leq \text{size\_}H_j$  loop  $\{ \text{End of } H_k : \text{append the remainder of } H_j \rightarrow H_k \}$ 
     $\langle \text{Idop}_j, \text{Sop}_j, \text{op}_j \rangle := H_j[i]$ ; 
    Append ( $H_k, \langle \text{Idop}_j, \text{Sop}_j, \text{op}_j \rangle$ );
     $i := i + 1$ ;
  endloop;

  {End of  $H_j$  : append the remainder of  $H_k \rightarrow H_j$ }
  while  $i \leq \text{size\_}H_k$  loop
     $\langle \text{Idop}_k, \text{Sop}_k, \text{op}_k \rangle := H_k[i]$ ; 
    Append ( $H_j, \langle \text{Idop}_k, \text{Sop}_k, \text{op}_k \rangle$ );
     $i := i + 1$ ;
  endloop;
end MOT2 ;

```

Example: Reconciliation **without using MOT2** (1)

(reconciliation without a master site)



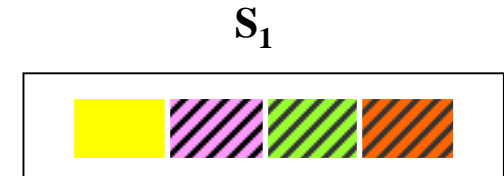
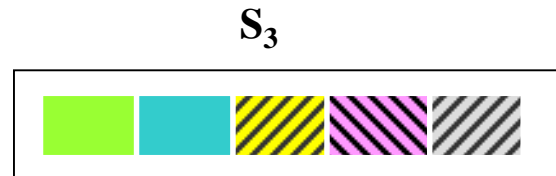
Reconciling copies of sites S_1 and S_3 ?

Example: Reconciliation **without using MOT2** (2)



Impossible to reconcile these histories !

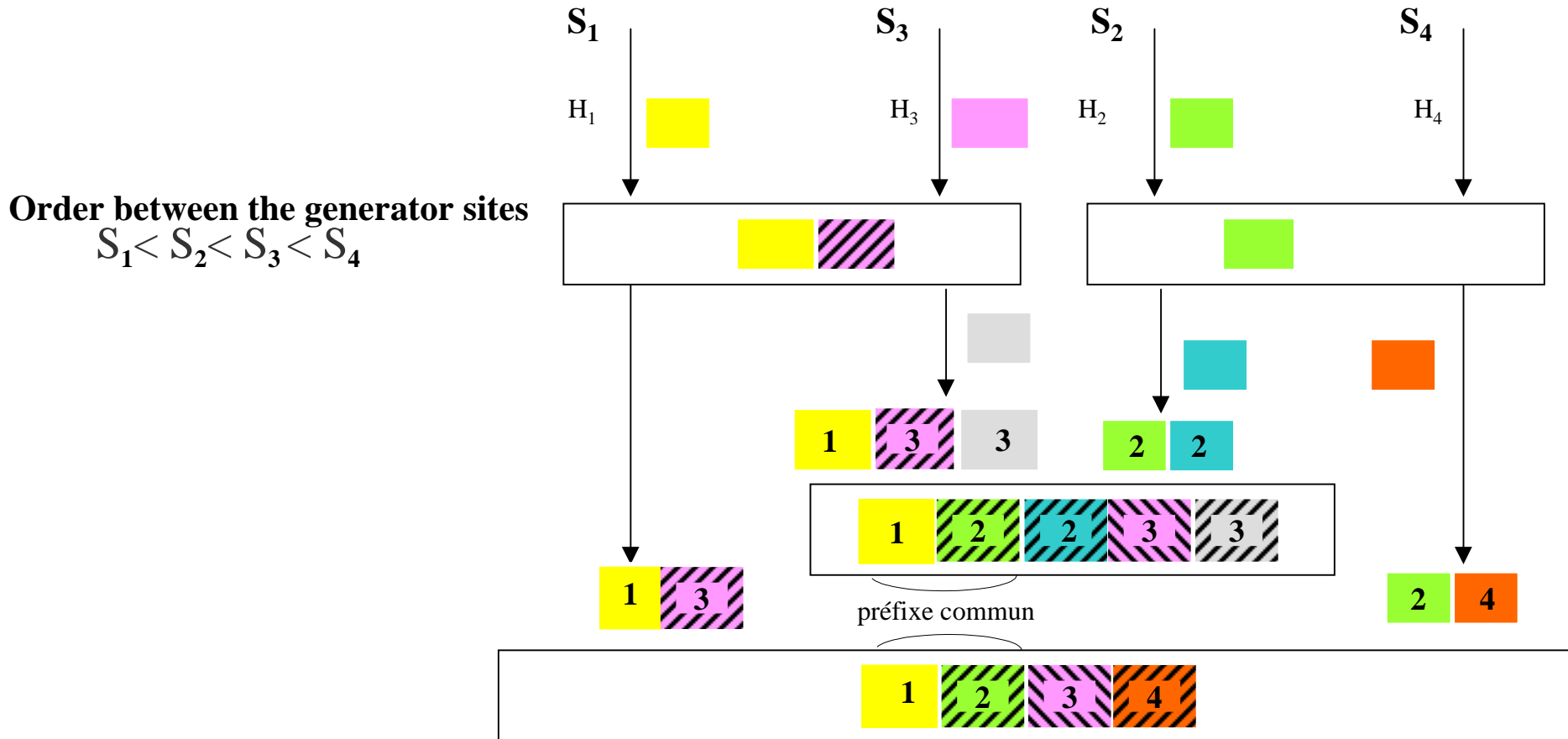
Histories to be merged



We cannot reconcile these histories because, although they have common operations ( ,  et ), they **do not appear**:

- in common prefix to both histories,
- in the same order,
- in the same form (because of forward transpositions).

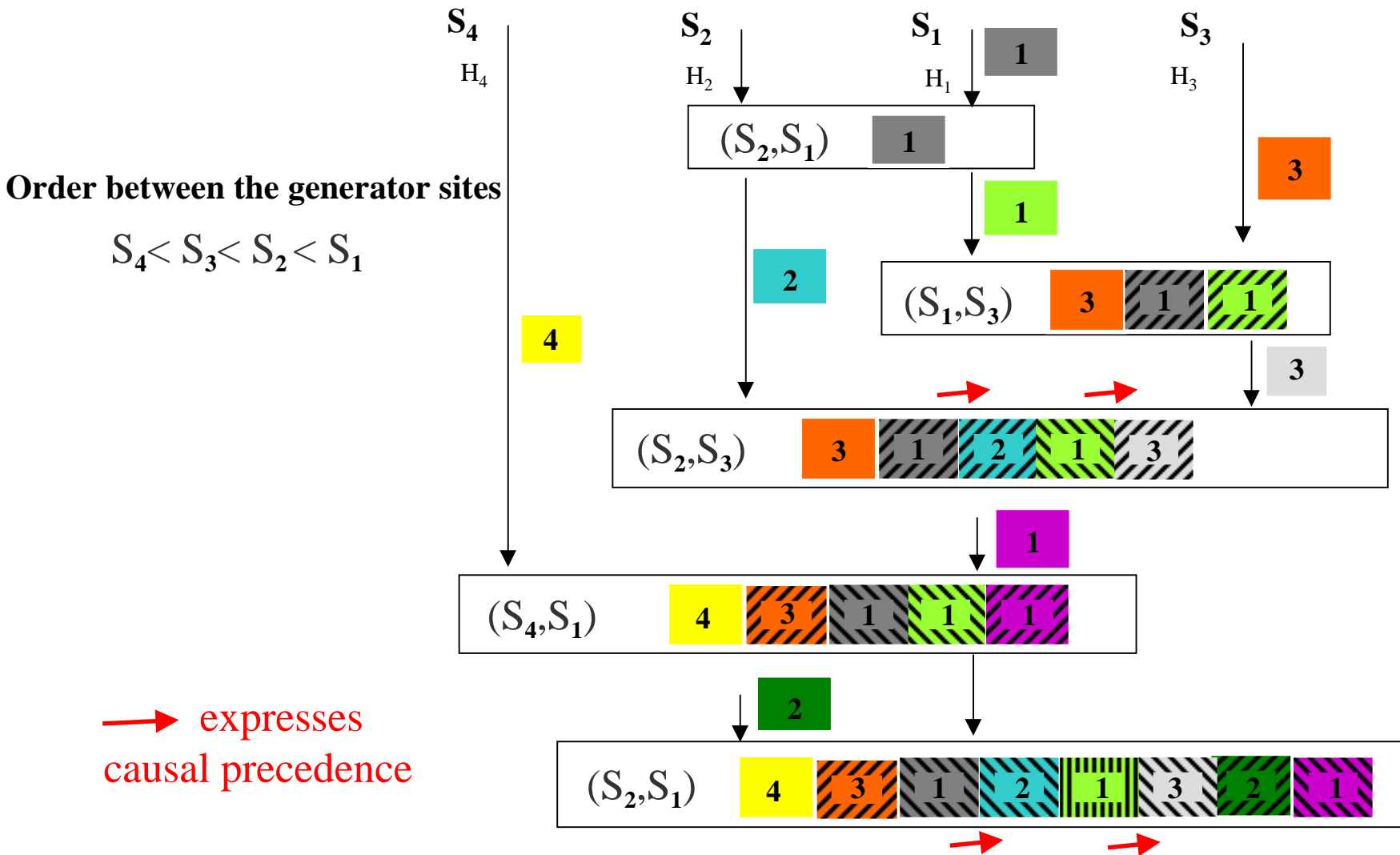
Example: Reconciliation by using MOT2 (3)



Reconciling S_1 (or S_4) and S_3 (or S_2) :



Example: Reconciliation by using MOT2 (4)



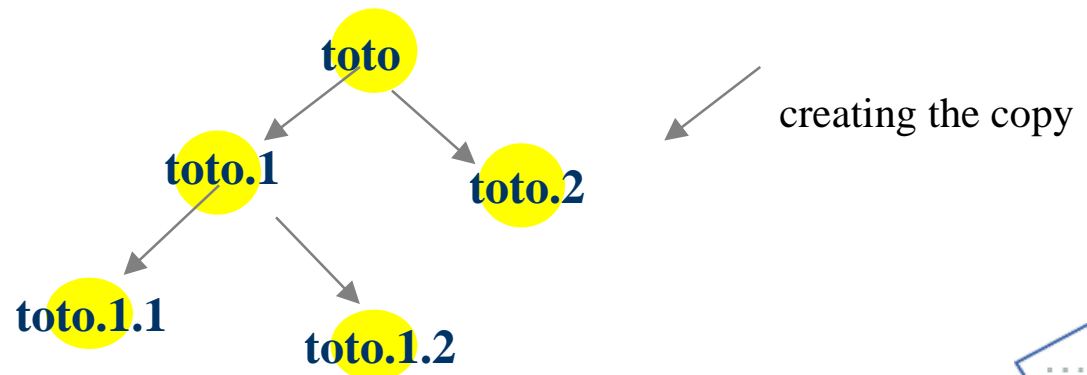
Ordering relation between sites / copies

Aim → Name of $\begin{bmatrix} \text{site} \\ \text{copy} \end{bmatrix}$ is unique and not re-allocated

How to obtain that ?

Site name = @IP → Problem in P2P when @IP allocated dynamically

Name of copy → Hierarchical identification schema



Conclusion

The MOT2 Merge Algorithm:

- reconciles pair-wise, any number of copies, while achieving convergence
 - enables free propagation of update operations
 - is based on Operational Transformation
 - constructs a unique global order between operations
(respecting the causal precedence relation)
without requiring any ordering mechanism
(state vectors, timestamps, sequencer, ..)
without any centralized component
without needing broadcast protocol
without needing consensus protocol
- scalable to a P2P collaborative environment