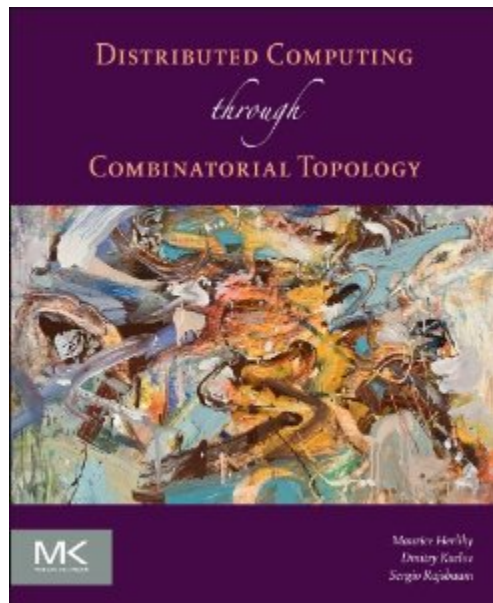


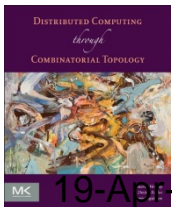
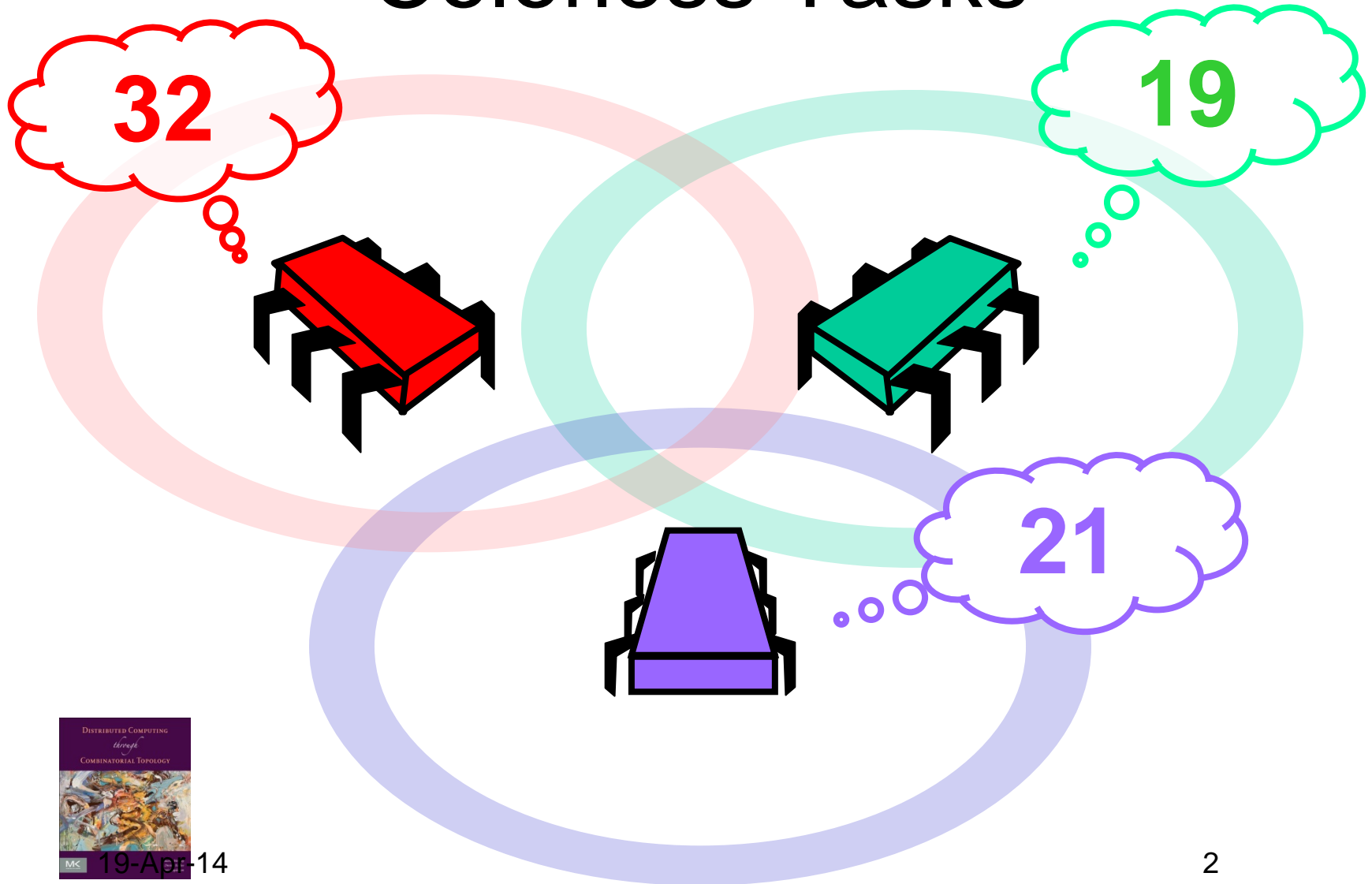
# Colorless Wait-Free Computation



Companion slides for  
**Distributed Computing  
Through Combinatorial Topology**  
Maurice Herlihy & Dmitry Kozlov & Sergio Rajsbaum  
Distributed Computing through  
Combinatorial Topology

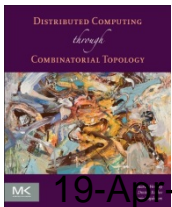
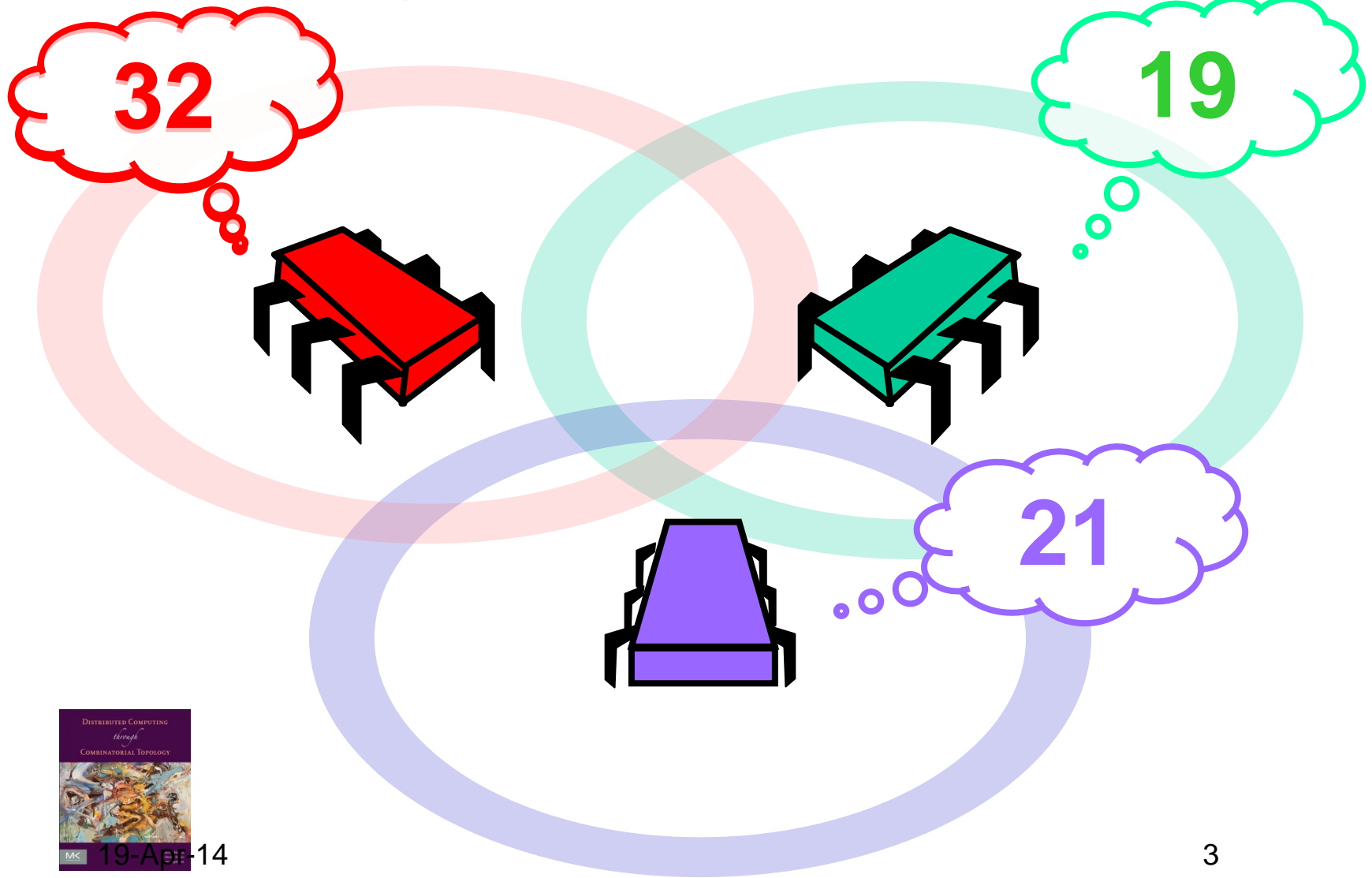


# Colorless Tasks





# Colorless Tasks



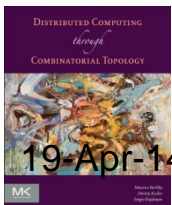
# Colorless Tasks

The *set* of input values ...

determines the *set* of output values.

Number and identities irrelevant...

for both input and output values

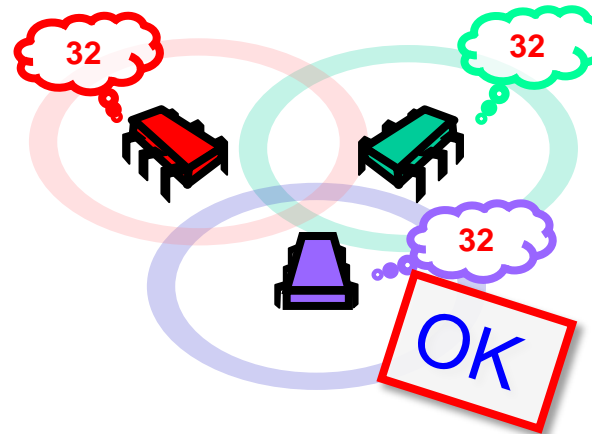


19-Apr-14

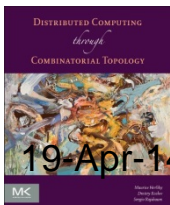
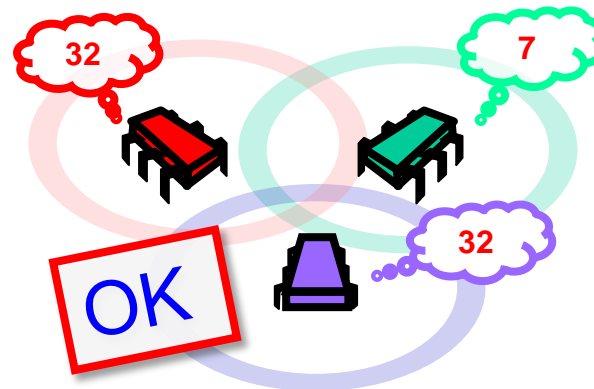


# Examples

Consensus



$k$ -set agreement

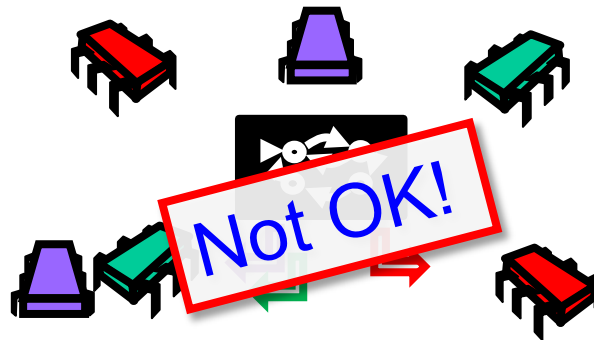


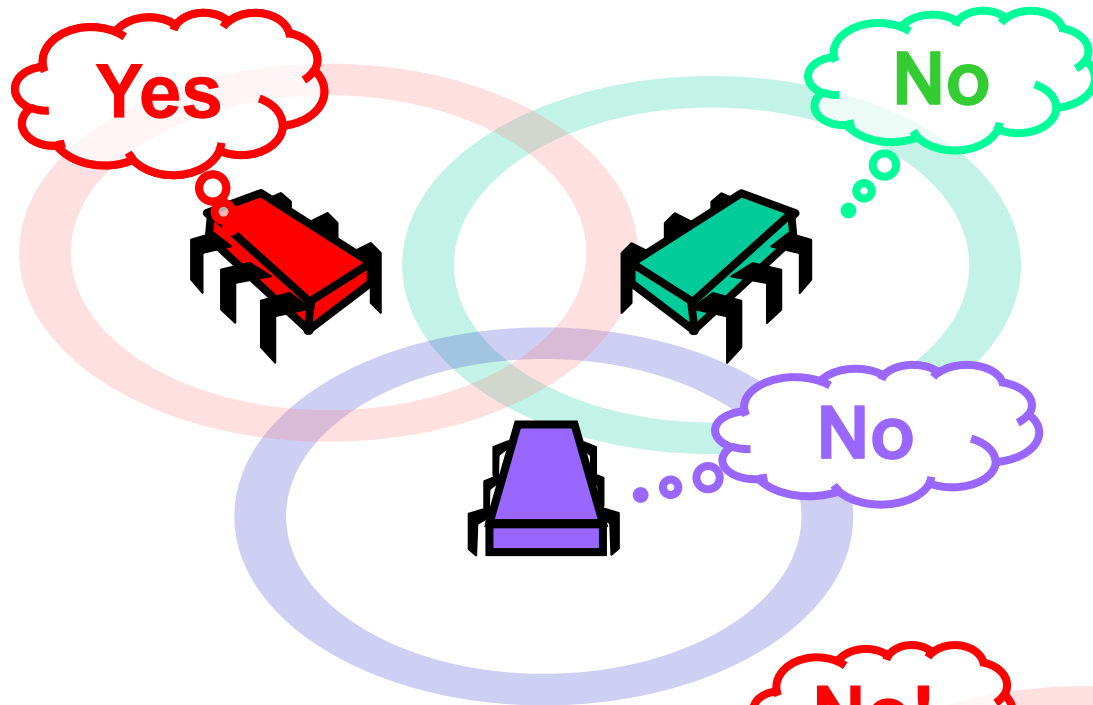
# Non-Examples

Weak Symmetry-Breaking

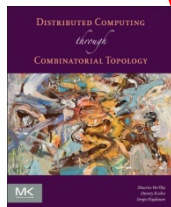
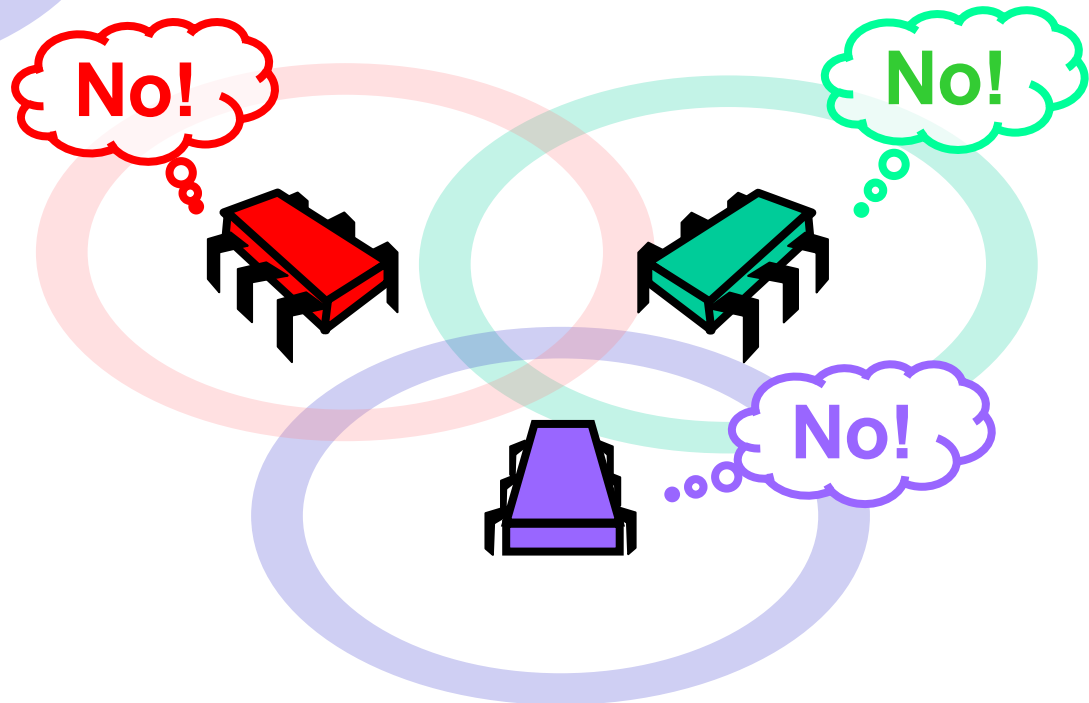
When all participate ...

At least one on group 0, group 1





Majority  
Not OK!

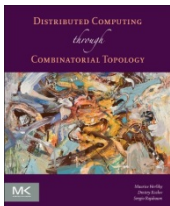


# Road Map

Operational Model

Combinatorial Model

Main Theorem

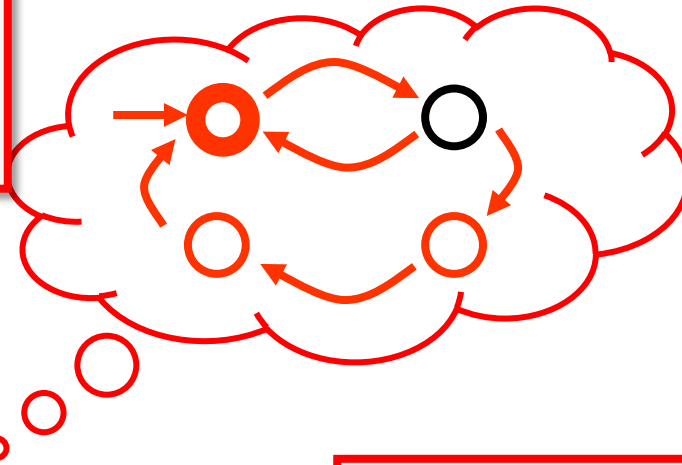


Distributed Computing through  
Combinatorial Topology

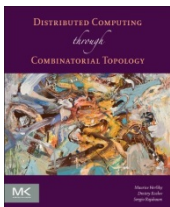
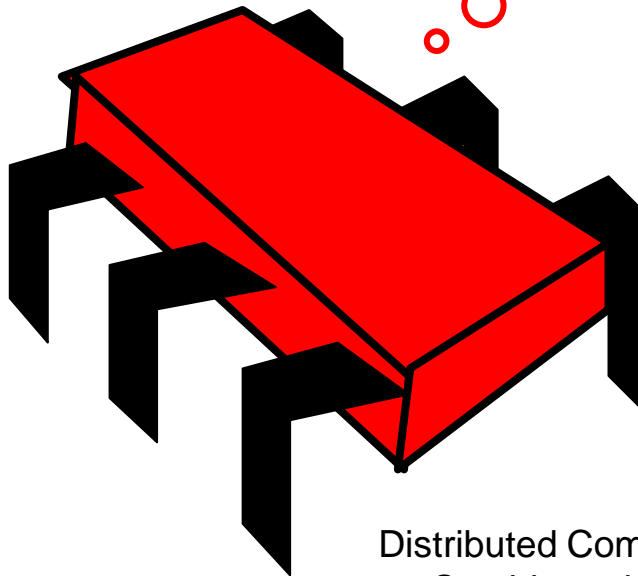


# Processes

A *process* is a state machine



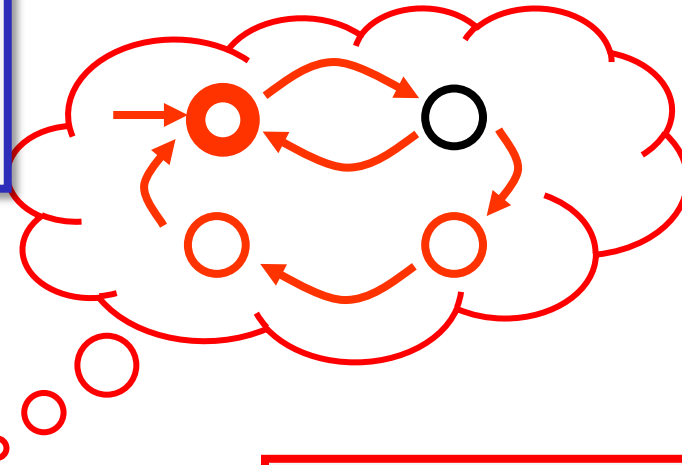
Could be Turing machines or more powerful



Distributed Computing through  
Combinatorial Topology

# Processes

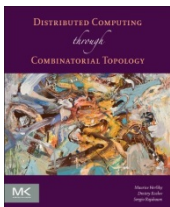
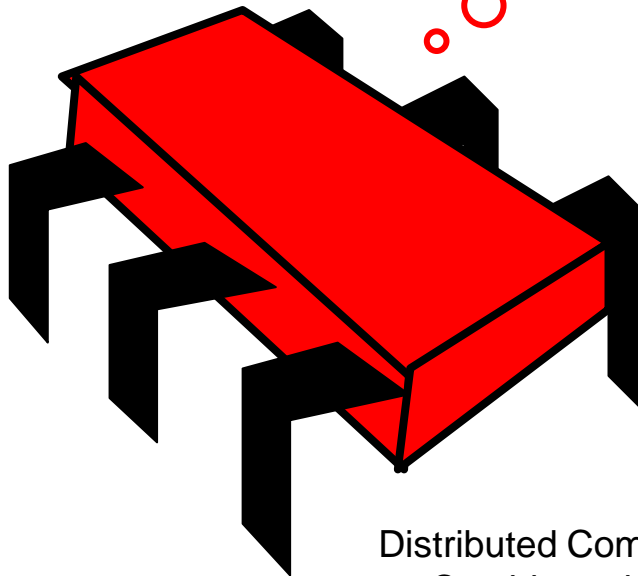
A process's state is called its *view*



Process names taken from a domain  $\Pi$

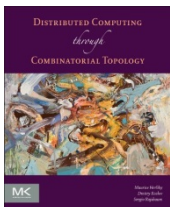
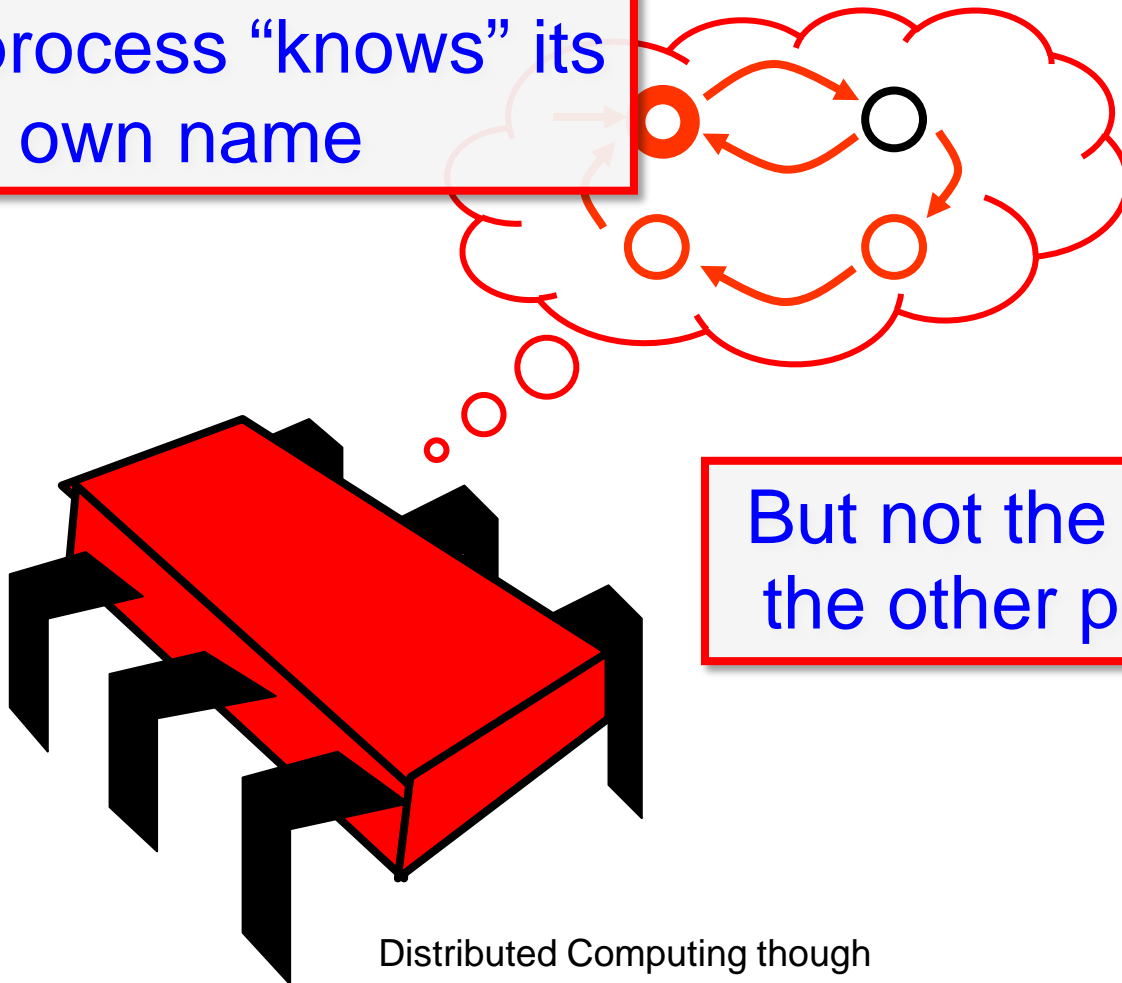
Each process has a unique *name* (color)

$$P_i \in \Pi$$



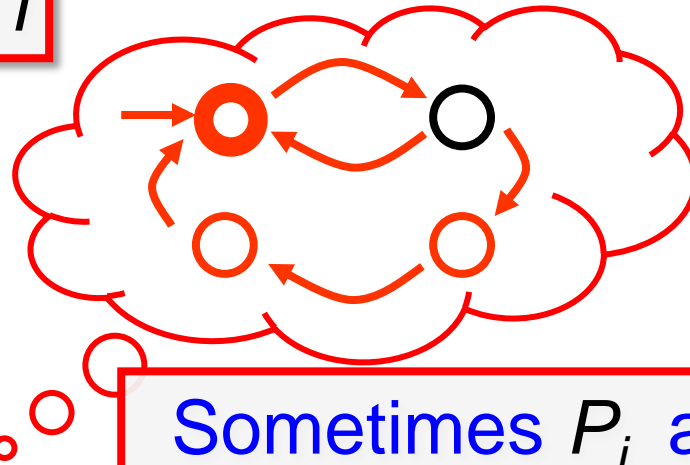
# Processes

Each process “knows” its  
own name

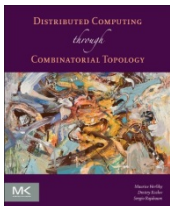
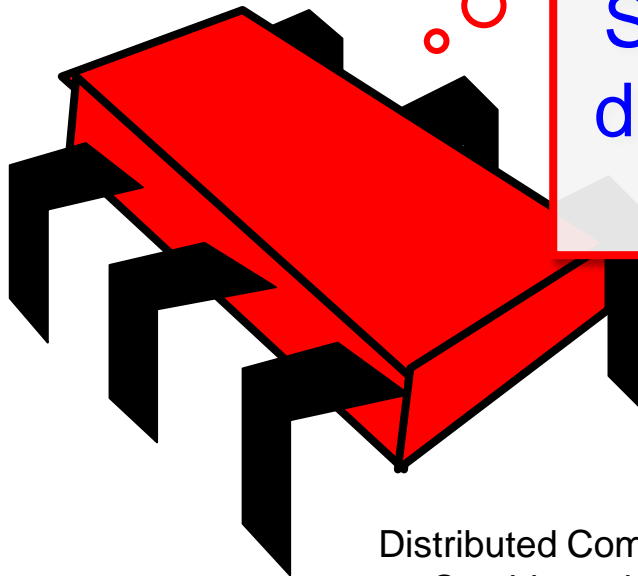


# Processes

Often,  $P_i$  is just  $i$

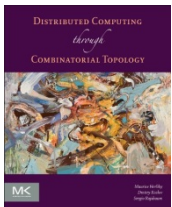
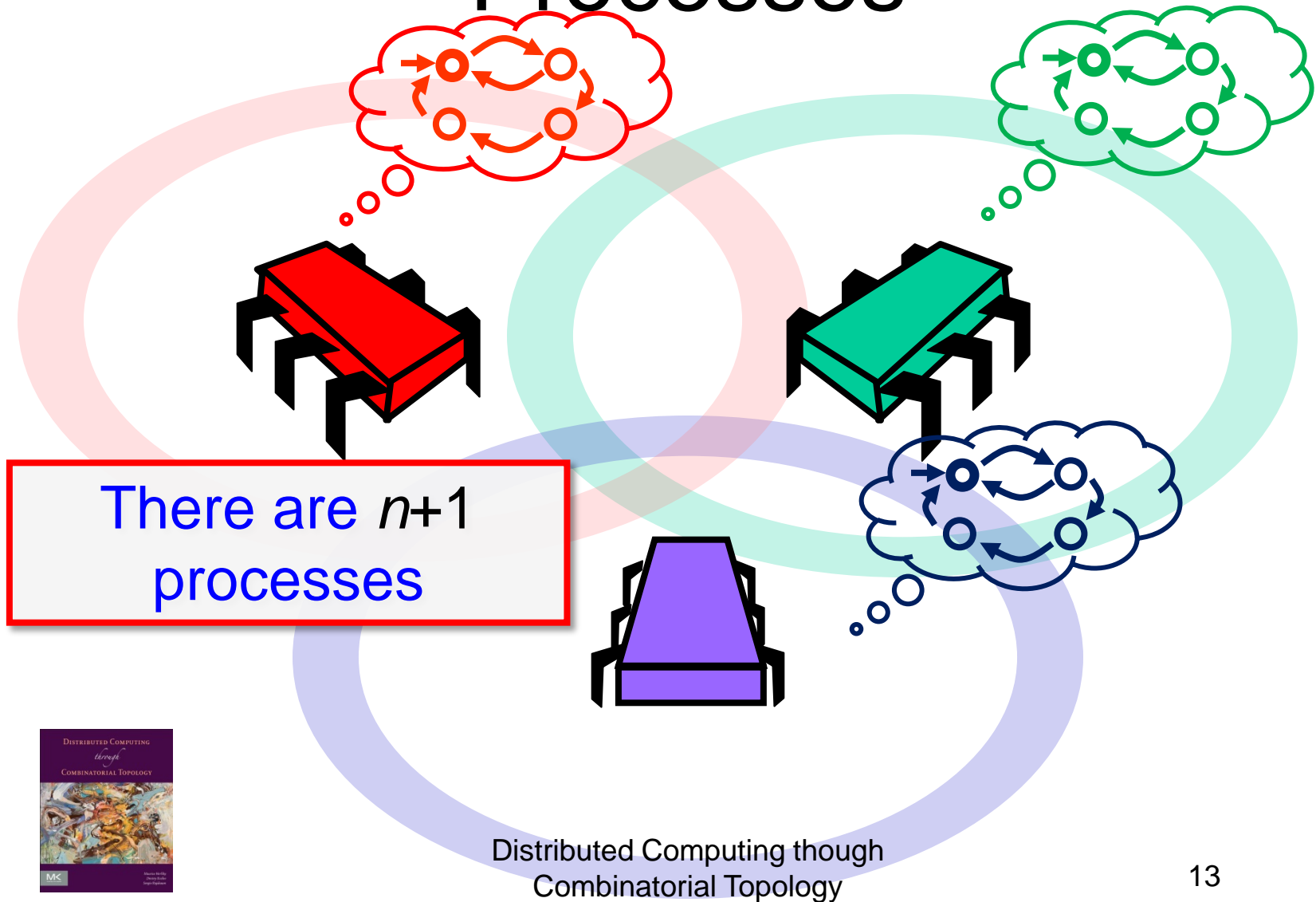


Sometimes  $P_i$  and  $i$  are distinct, and the process “knows”  $P_i$  but not  $i$

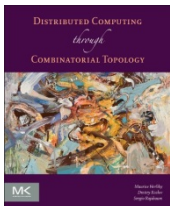




# Processes



# Shared Read-Write Memory

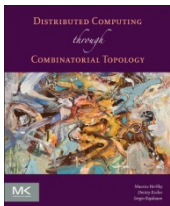


# Immediate Snapshot

Individual reads & writes are too low-level ...

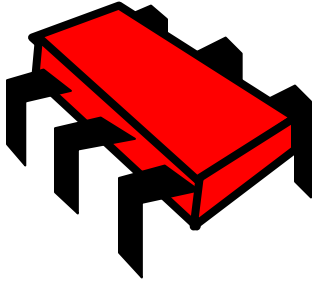
A *snapshot* = atomic read of all memory

We will use *immediate snapshot* ...





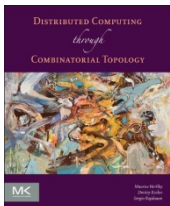
# Immediate Snapshot



write view to memory

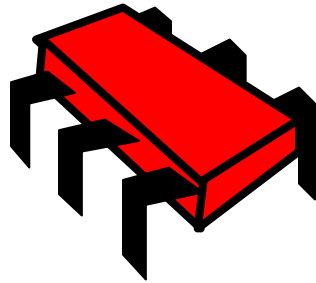
take snapshot

} adjacent steps!

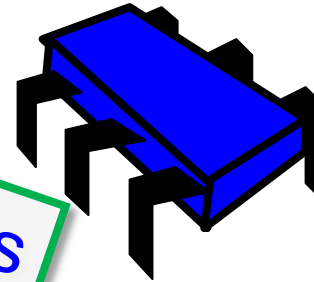




# Immediate Snapshot



Concurrent steps

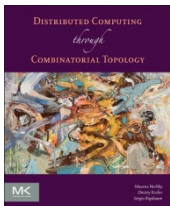


write view to memory

take snapshot

write view to memory

take snapshot



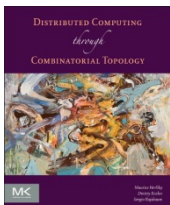


# Immediate Snapshot

```
immediate
```

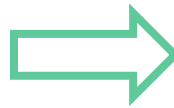
```
  mem[i] := view;
```

```
  snap := snapshot(mem[*])
```





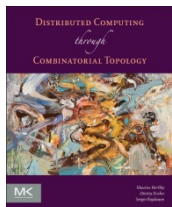
P	Q	R
write		
snap		
	write	
	snap	
		write
		snap
{p}	{p,q}	{p,q,r}



P	Q	R
write		
snap		
	write	write
	snap	snap
{p}	{p,q,r}	{p,q,r}



P	Q	R
write	write	write
snap	snap	snap
{p,q,r}	{p,q,r}	{p,q,r}





# Realistic?

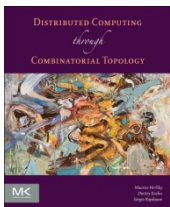
No!

My laptop reads only a few contiguous memory words at a time

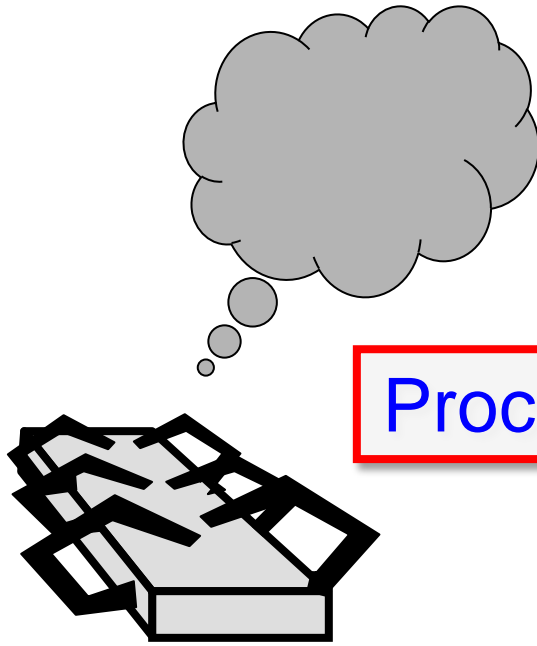
Yes!

Simpler lower bounds: if it's impossible with IS, it's impossible on your laptop.

Can implement IS from read-write



# Crashes

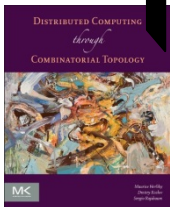
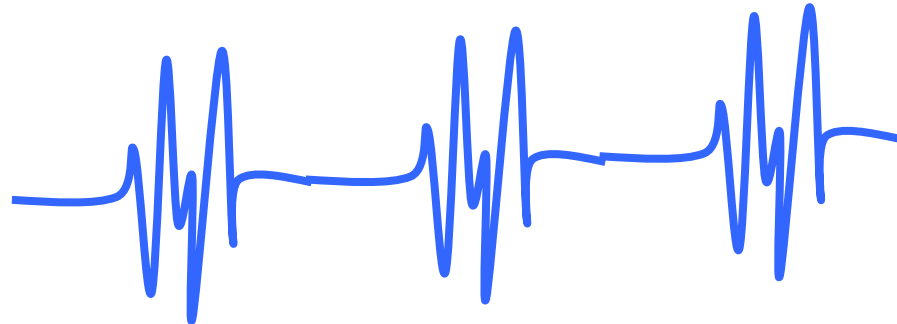
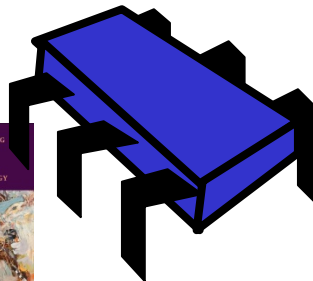
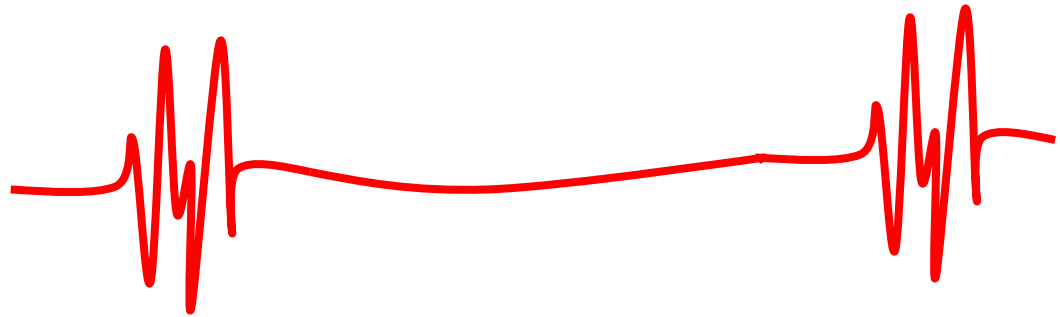
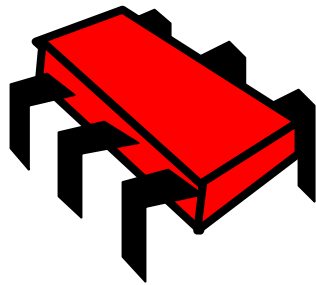
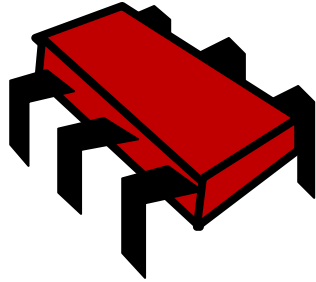


Processes may halt without warning

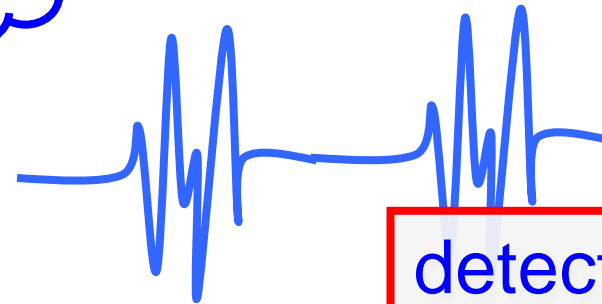
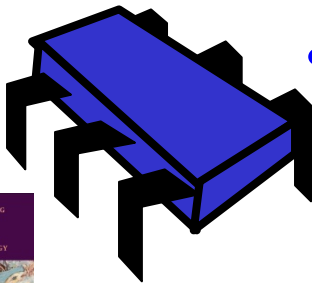
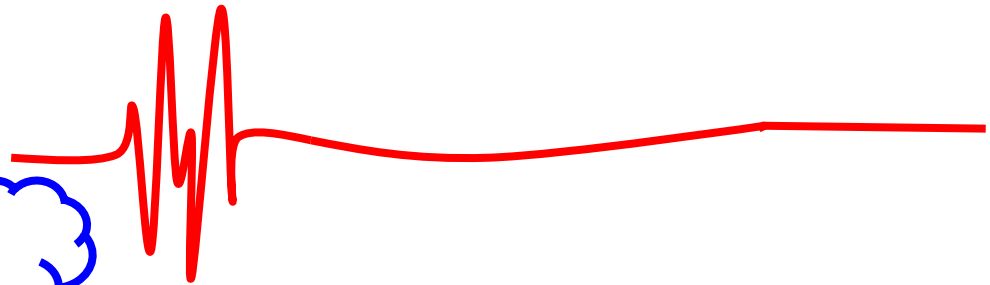
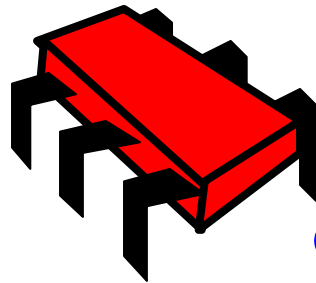
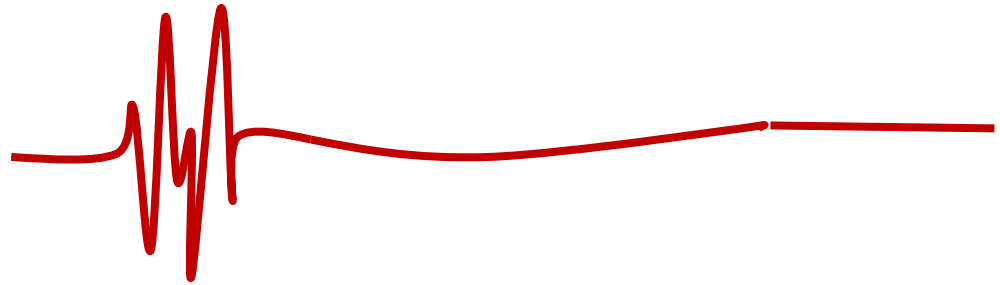
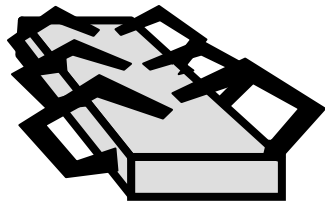
as many as  $n$  out of  $n+1$



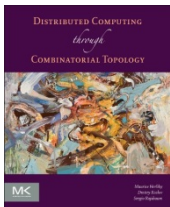
# Asynchronous



# Asynchronous Failures



detection impossible



# Configurations

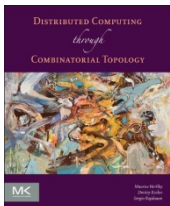
$$C = \{s_0, \dots, s_n\}$$



set of *simultaneous* process states

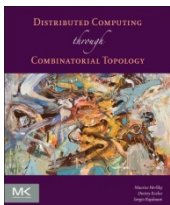
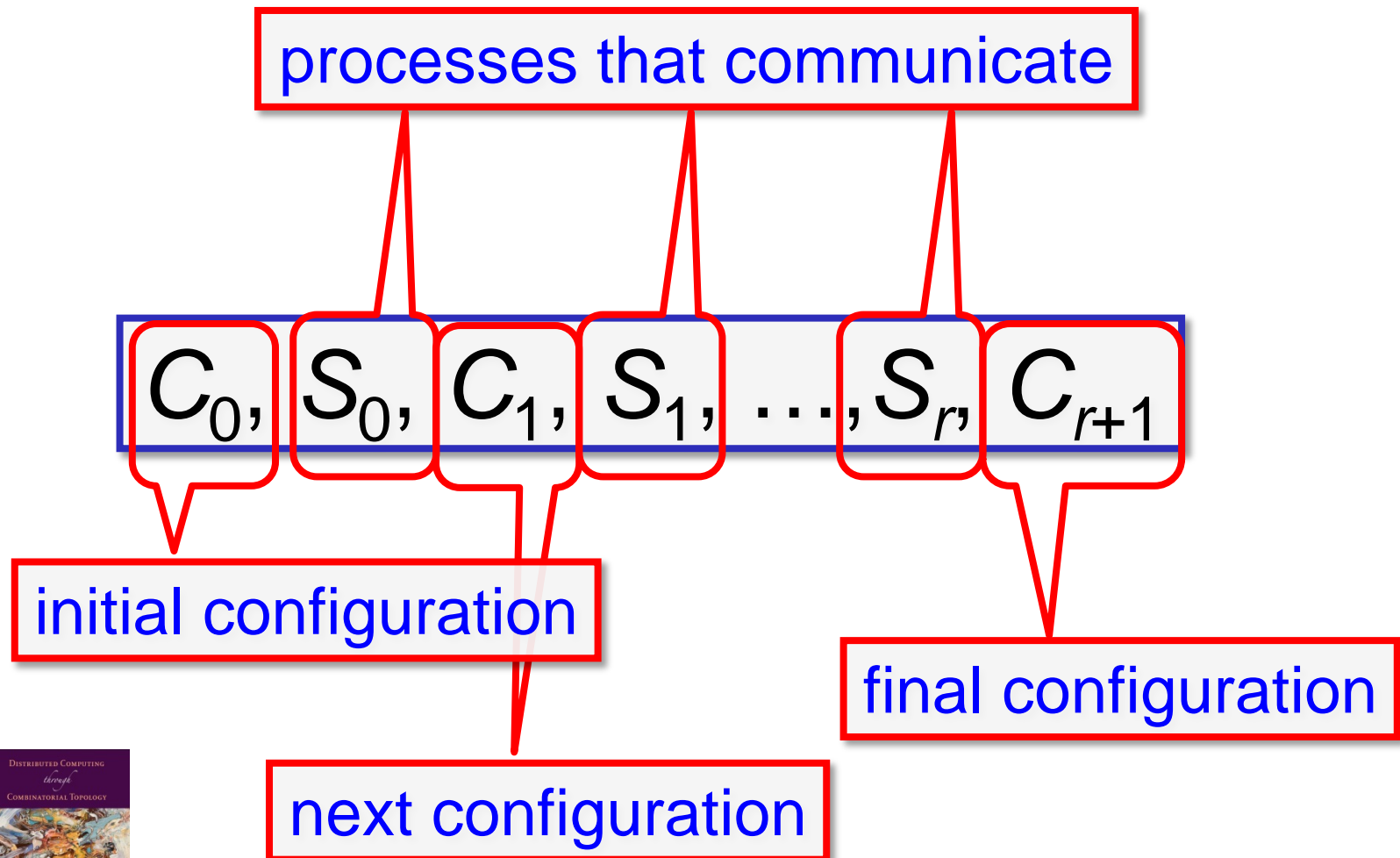
*initial* configurations

*final* configurations





# Executions



# Executions

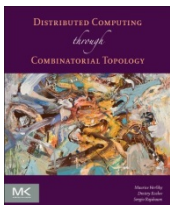
triple is a *concurrent step*

$C_0, S_0, C_1, S_1, \dots, S_r, C_{r+1}$

Processes in  $S_0$  said to *participate* in step

Only  $P_i \in S_0$  can change between  $C_0$  and  $C_1$

state change result of communication



# Executions

*finite*

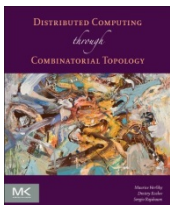
$C_0, S_0, C_1, S_1, \dots, S_r, C_{r+1}$

*infinite*

$C_0, S_0, C_1, S_1, \dots,$

*partial*

$C_0, S_0, C_1, S_1, \dots, S_r, C_{r+1}$

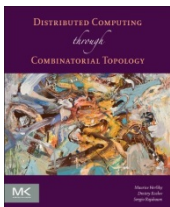


# Crashes are Implicit

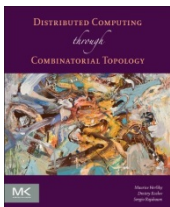
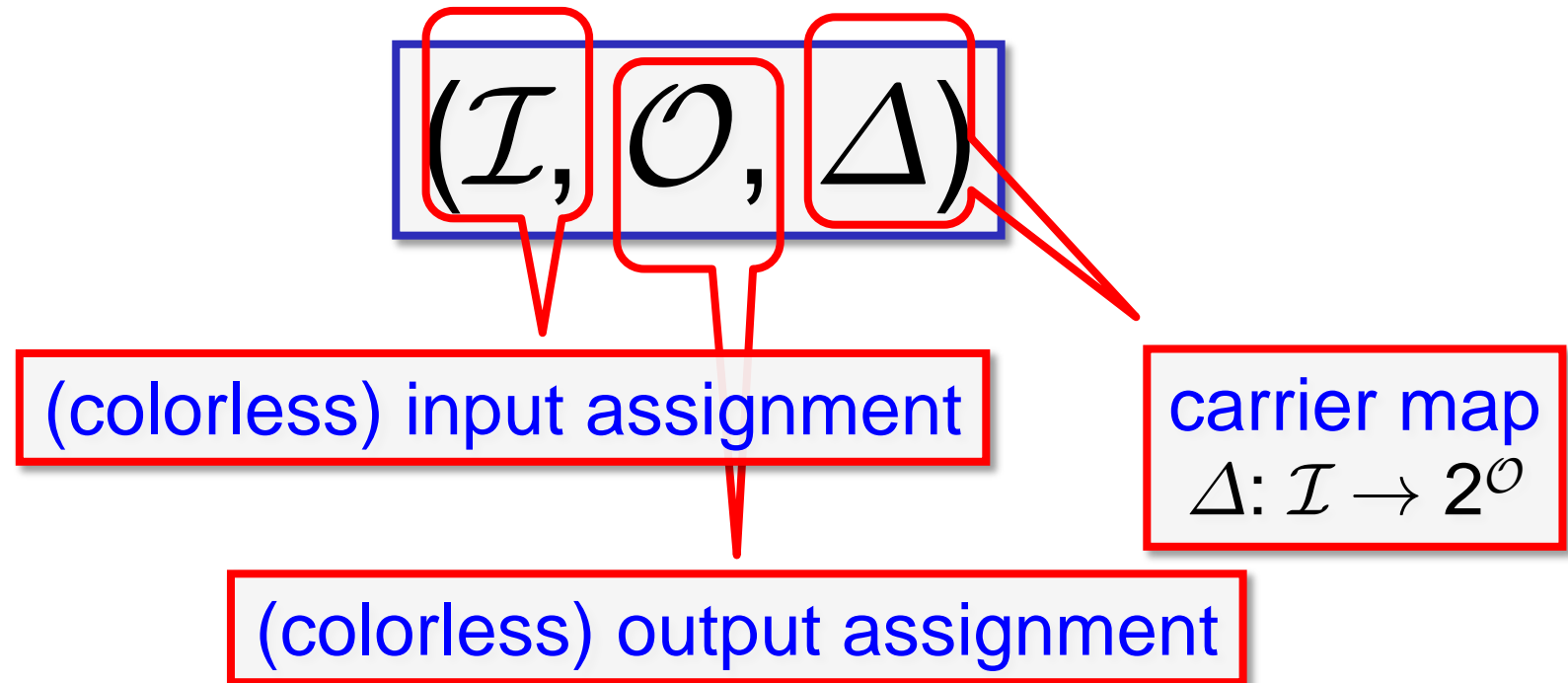
$C_0, S_0, C_1, S_1, \dots, S_r, C_{r+1}$

If  $P_i$ 's state not final in the final configuration  
then  $P_i$  has *crashed*.

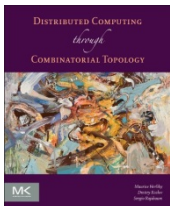
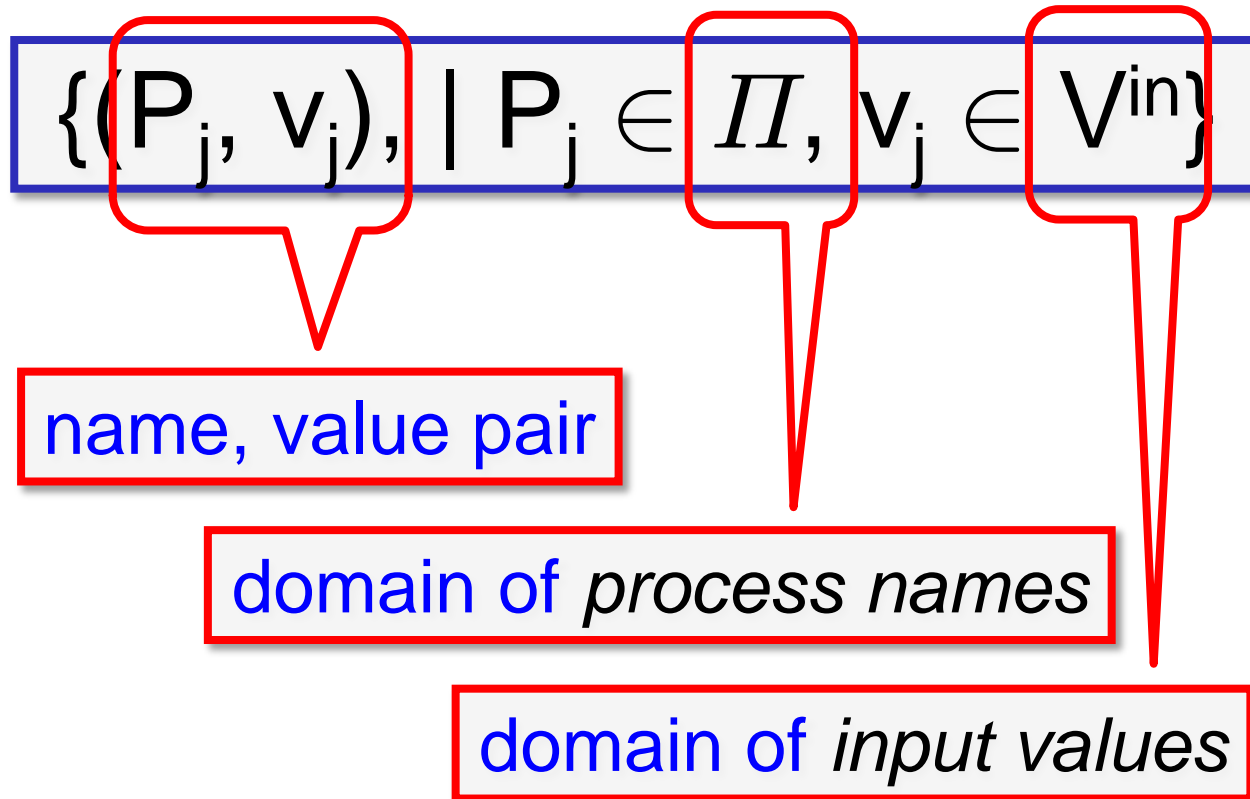
crash cannot be detected in finite time



# Colorless Tasks



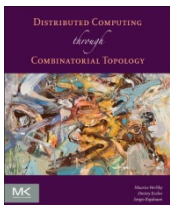
# Input Assignments



# Colorless Input Assignments

$$\{(P_j, v_j), \mid P_j \in \Pi, v_j \in V^{\text{in}}\}$$

discard process names, keep values

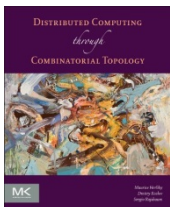




# (Colorless) Output Assignments

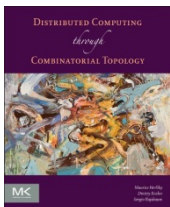
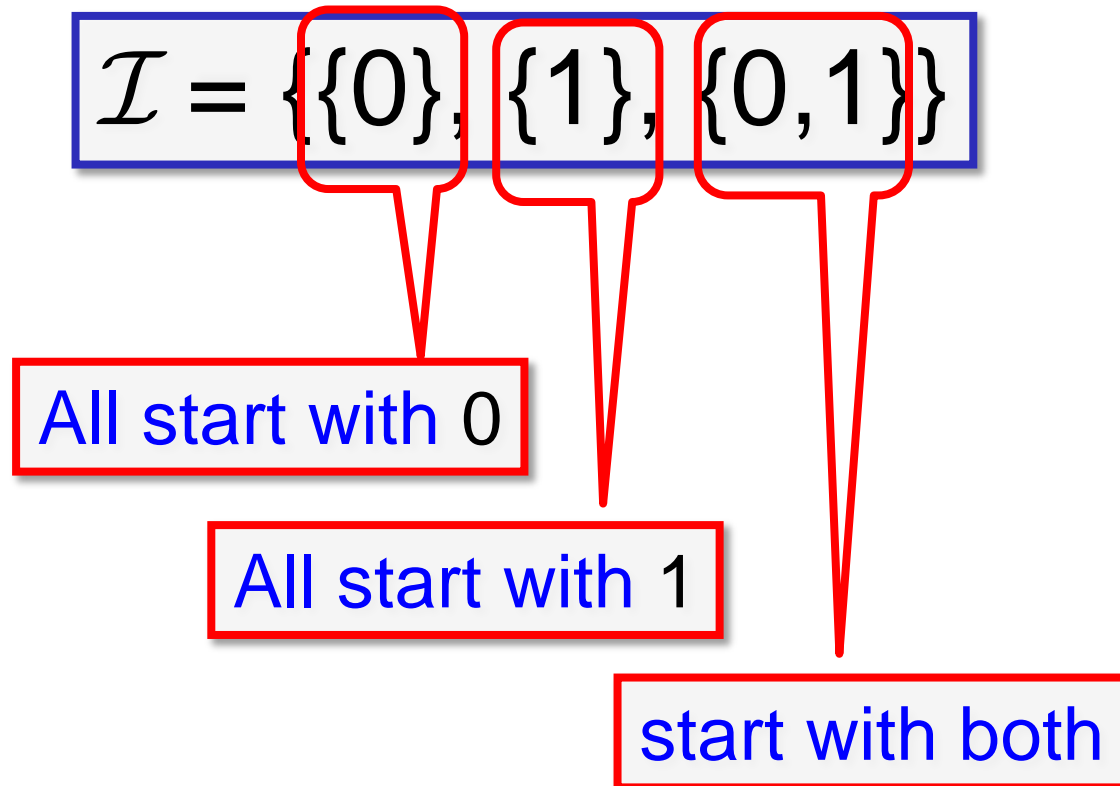
$$\{(P_j, v_j), \mid P_j \in \Pi, v_j \in V^{\text{out}}\}$$

$$\{(\textcolor{gray}{P}_j, \textcolor{gray}{v}_j), \mid \textcolor{gray}{P}_j \in \textcolor{gray}{\Pi}, v_j \in V^{\text{out}}\}$$





# Example: Binary Consensus



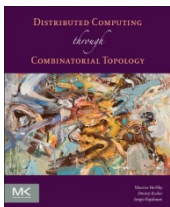
# Example: Binary Consensus

$$\mathcal{I} = \{\{0\}, \{1\}, \{0,1\}\}$$

$$\mathcal{O} = \{\{0\}, \{1\}\}$$

All decide 0

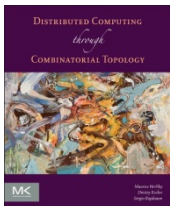
All decide 1



# Example: Binary Consensus

$$\Delta(\{0\}) = \{\{0\}\}$$

All start with 0,  
all decide 0

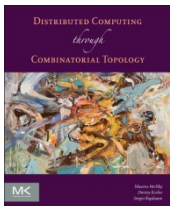


# Example: Binary Consensus

$$\Delta(\{0\}) = \{\{0\}\}$$

$$\Delta(\{1\}) = \{\{1\}\}$$

All start with 1,  
all decide 1



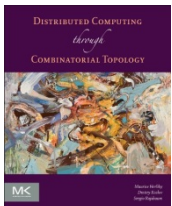
# Example: Binary Consensus

$$\Delta(\{0\}) = \{\{0\}\}$$

$$\Delta(\{1\}) = \{\{1\}\}$$

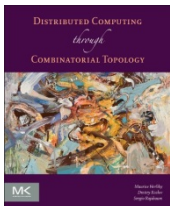
$$\Delta(\{0,1\}) = \{\{0\},\{1\}\}$$

with mixed inputs,  
all decide 0,  
or all decide 1



# Colorless Layered Protocol

```
shared mem array 0..N-1, 0..n of Value
view := input
for l := 0 to N-1 do
  immediate
    mem[l][i] := view;
    snap := snapshot(mem[l][*])
  view := set of values in snap
return  $\delta$ (view)
```



# Colorless Layered Protocol

shared **mem** array  $0..N-1, 0..n$  of **Value**

`view := input`

`for j := 0 to  $n-1$`

`immediate`

`mem[j][i]`

`snap := snapshot(mem[j][0..n-1])`

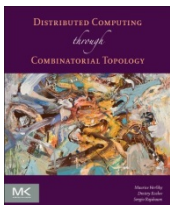
`view := set`

`return  $\delta$ (view)`

*2-dimensional memory array*

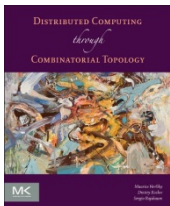
*row is clean per-layer memory*

*column is per-process word*



# Colorless Layered Protocol

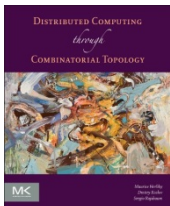
```
shared mem array 0..N-1, 0..n of Value
view := input
for j := 0 to N-1 do
  immediate
  initial view is input value
  snap := snapshot(mem[j][*])
  view := set of values in snap
return  $\delta$ (view)
```





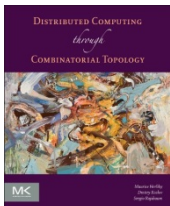
# Colorless Layered Protocol

```
shared mem array 0..N-1, 0..n of Value
view := input
for  $l := 0$  to  $N-1$  do
  immediate
  mem[ $i$ ][ $i$ ] := view;
  snap run for  $N$  layers mem[ $j$ ][*])
  view := set of values in snap
return  $\delta$ (view)
```



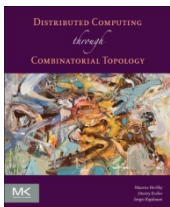
# Colorless Layered Protocol

```
shared
view :
for j := 0 to N-1 do
  immediate
    mem[l][i] := view;
    snap := snapshot(mem[l][*])
  view := set of values in snap
return  $\delta$ (view)
```



# Colorless Layered Protocol

```
shared mem array 0..N-1, 0..n of Value
view := input
for j := 0 to N-1 do
    immediate new view is set of values seen
    mem[j][i] := view;
    snap := snapshot(mem[j][*])
view := set of values in snap
return  $\delta$ (view)
```



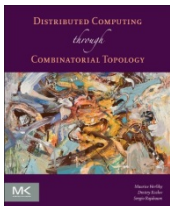
# Colorless Layered Protocol

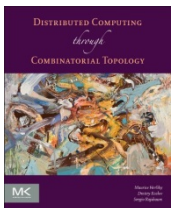
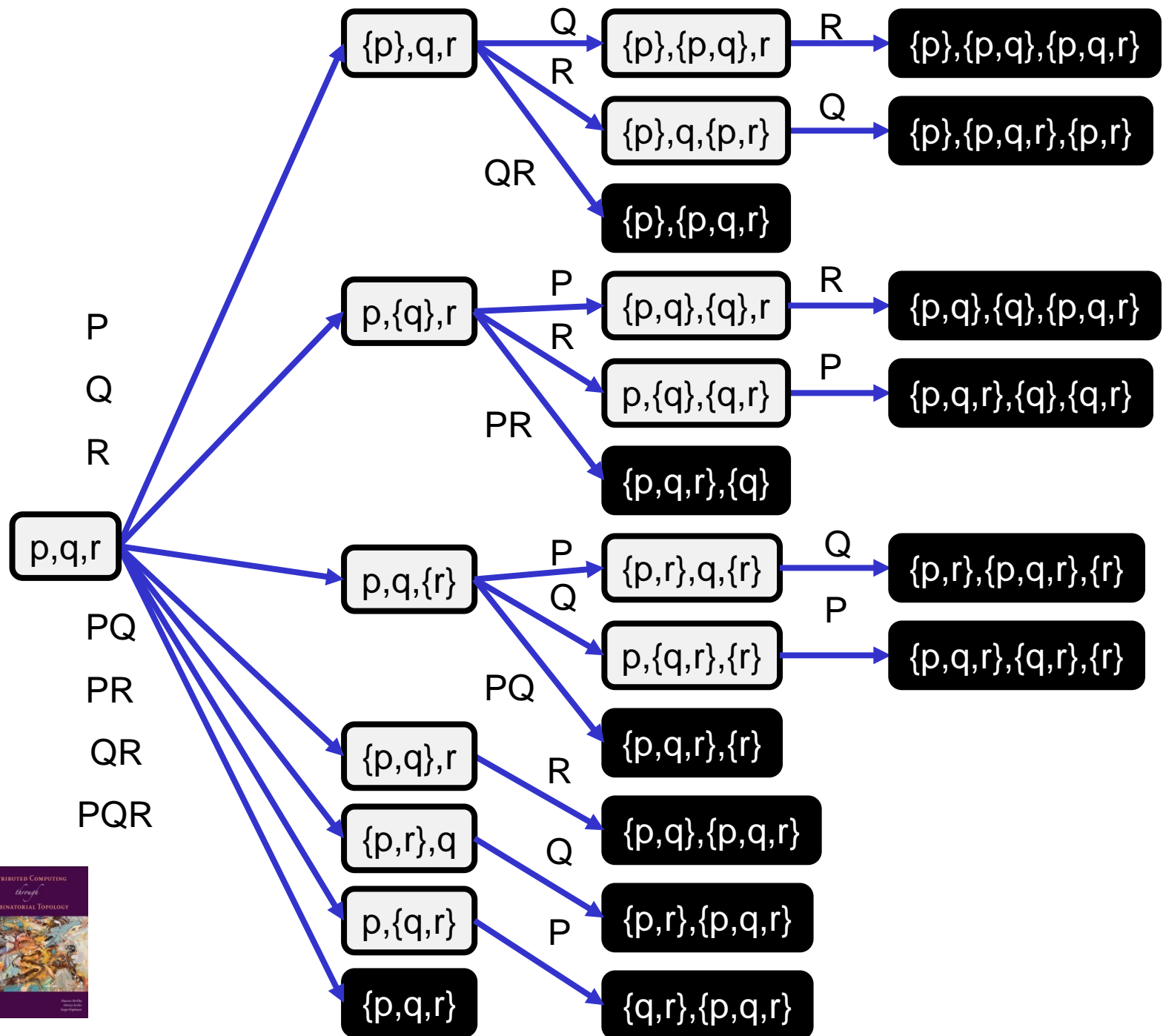
```
shared mem array 0..N-1, 0..n of Value  
view := input  
for j := 0 to N-1 do  
  immediate
```

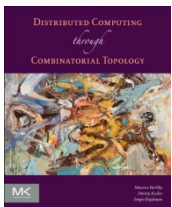
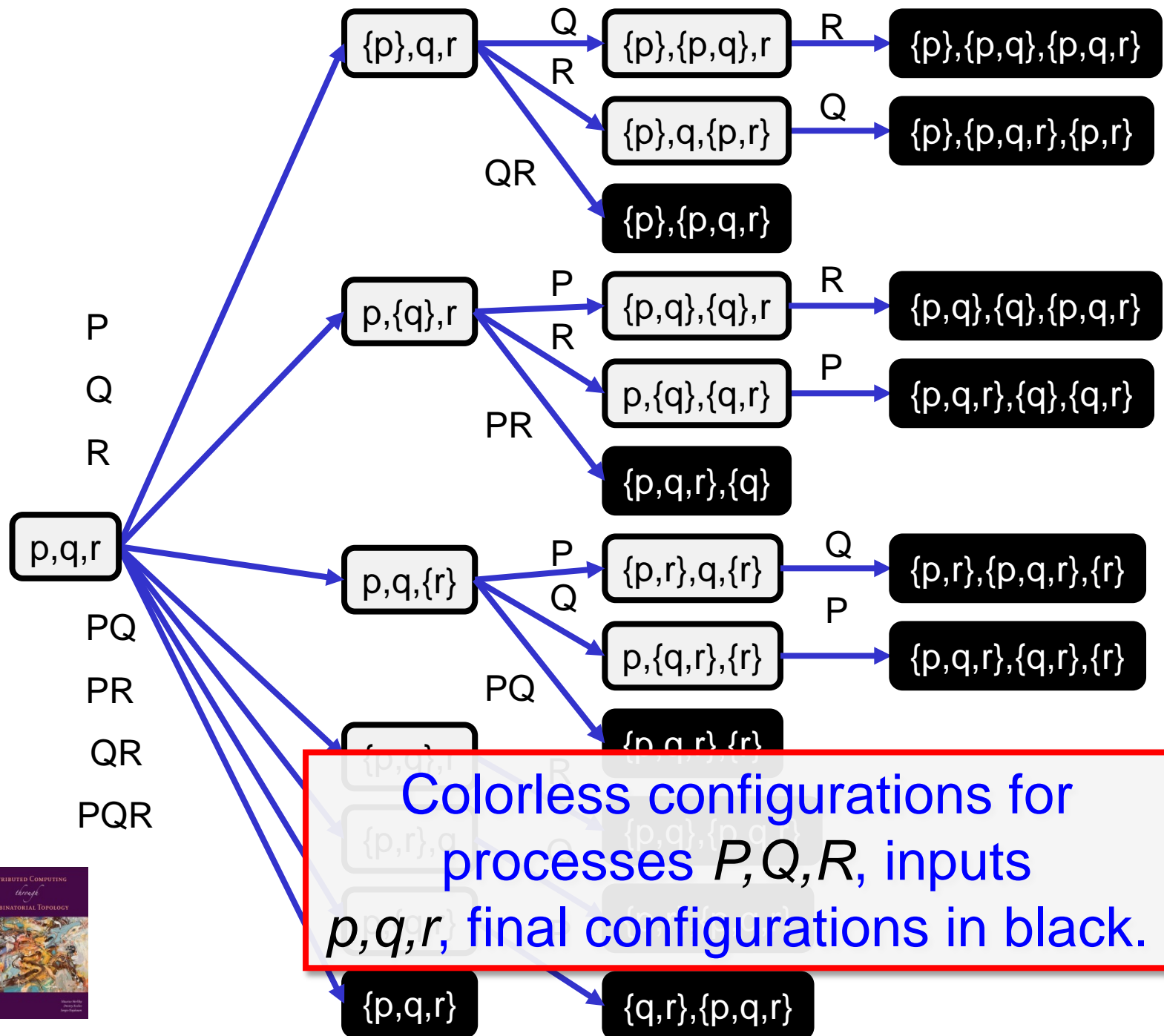
finally apply decision value to final view

```
  snap := snapshot(mem[j][*])  
  view := set of values in snap
```

```
return  $\delta$ (view)
```





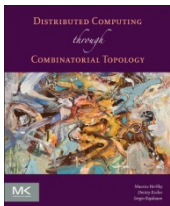


# Road Map

Operational Model

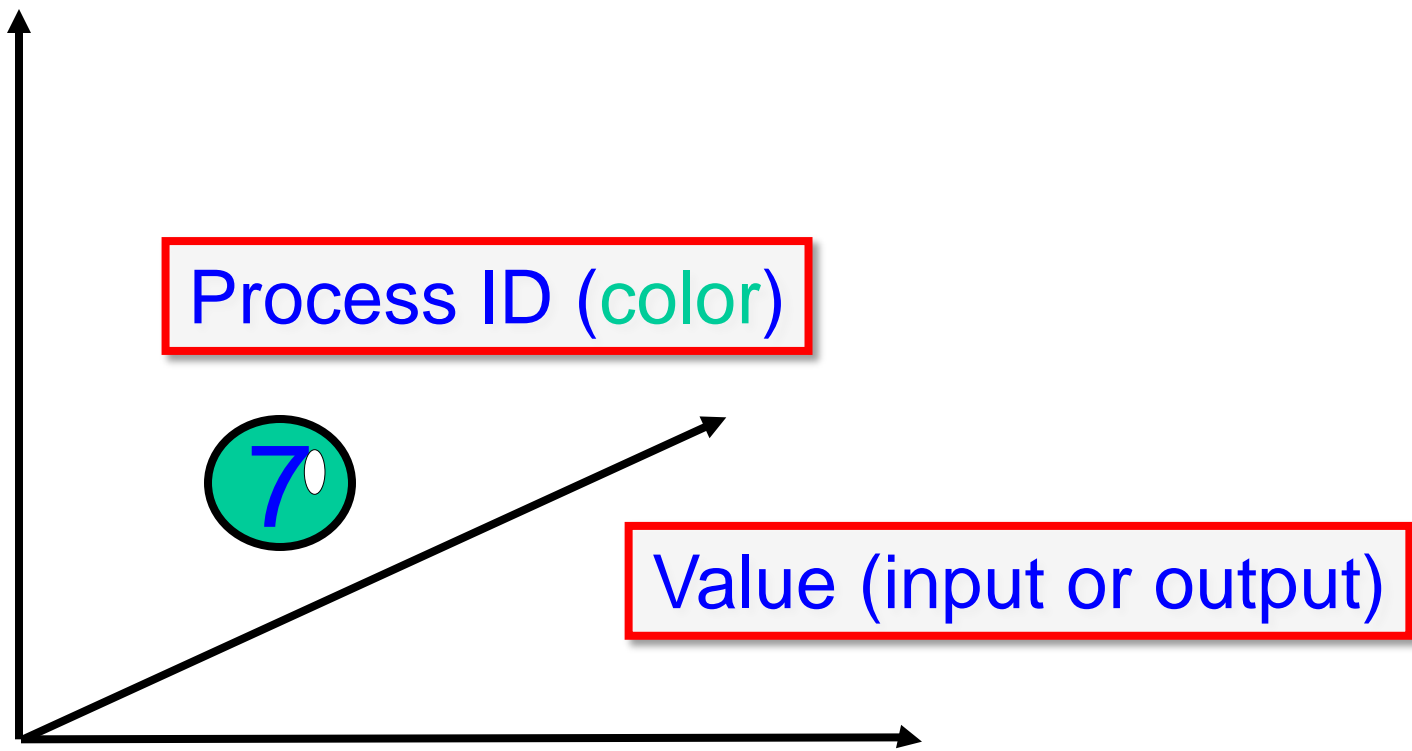
Combinatorial Model

Main Theorem



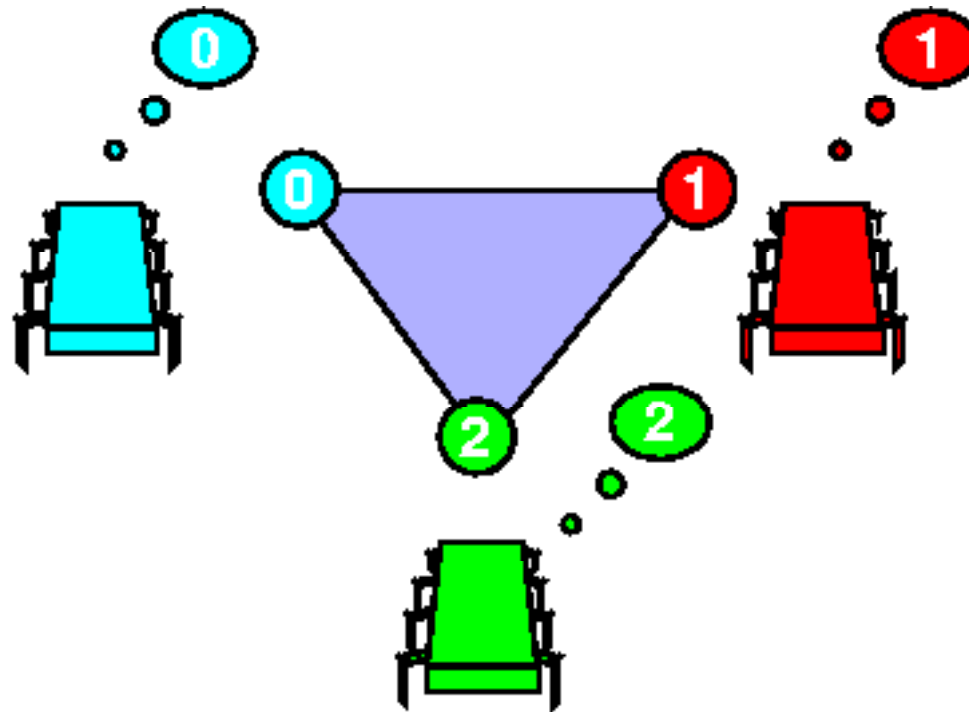
Distributed Computing through  
Combinatorial Topology

# Vertex = Process State

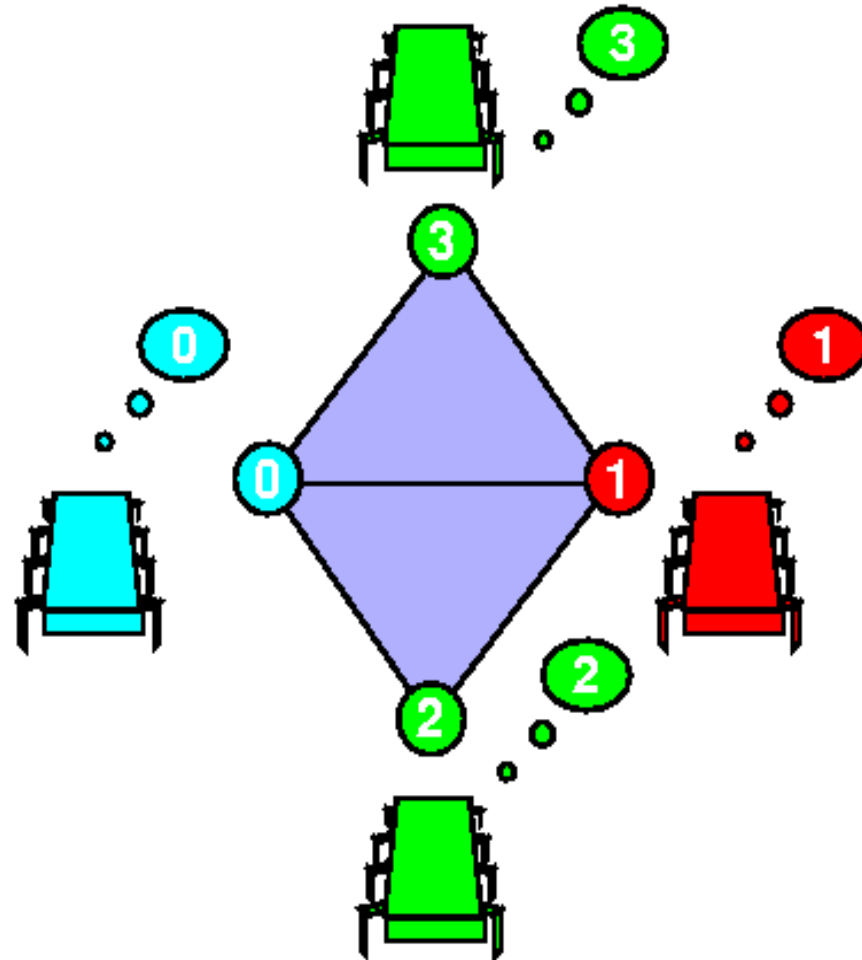




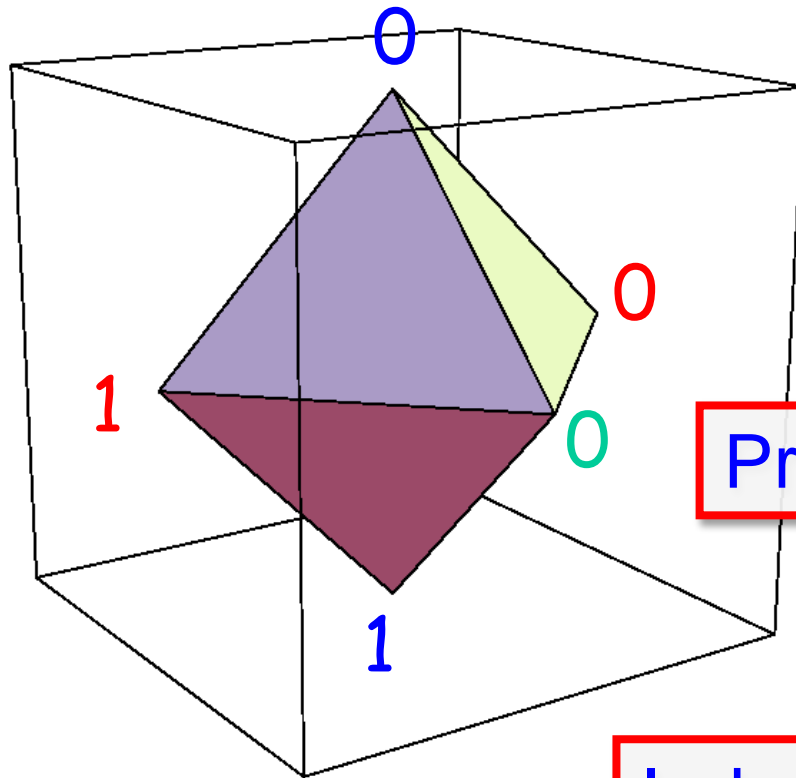
# Simplex = Global State



# Complex = Global States



# Input Complex for Binary Consensus

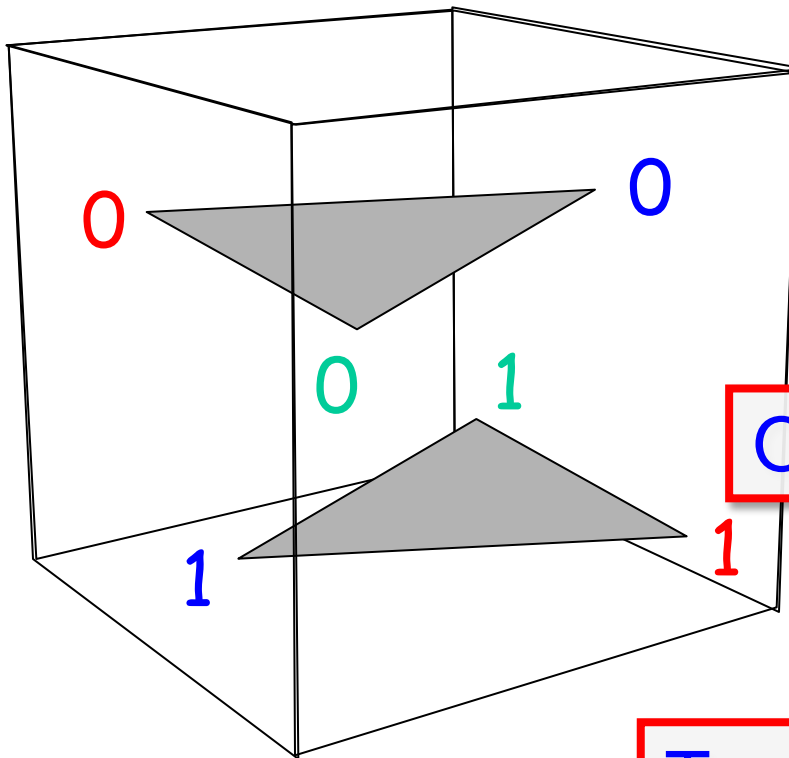


All possible initial states

Processes: red, green, blue

Independently assigned 0 or 1

# Output Complex for Binary Consensus

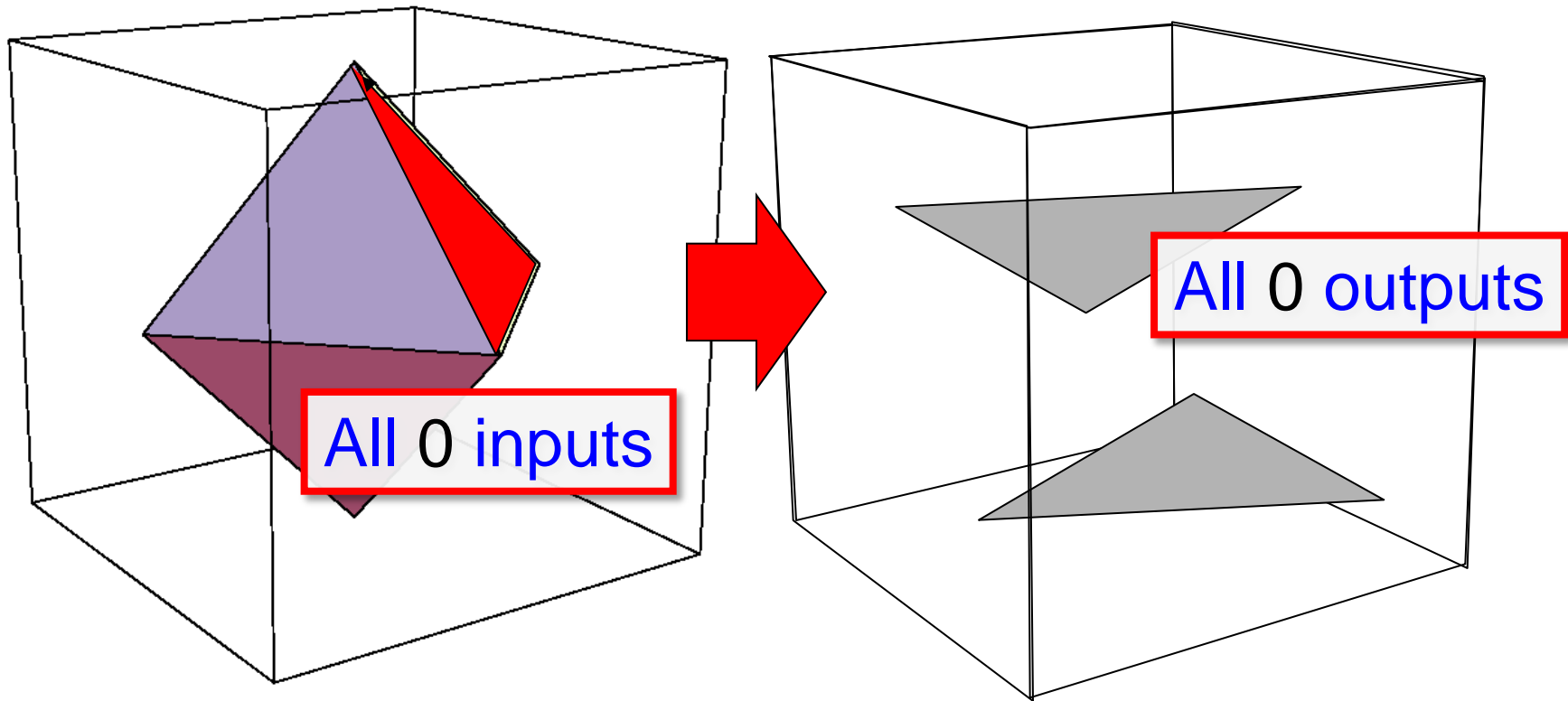


All possible final states

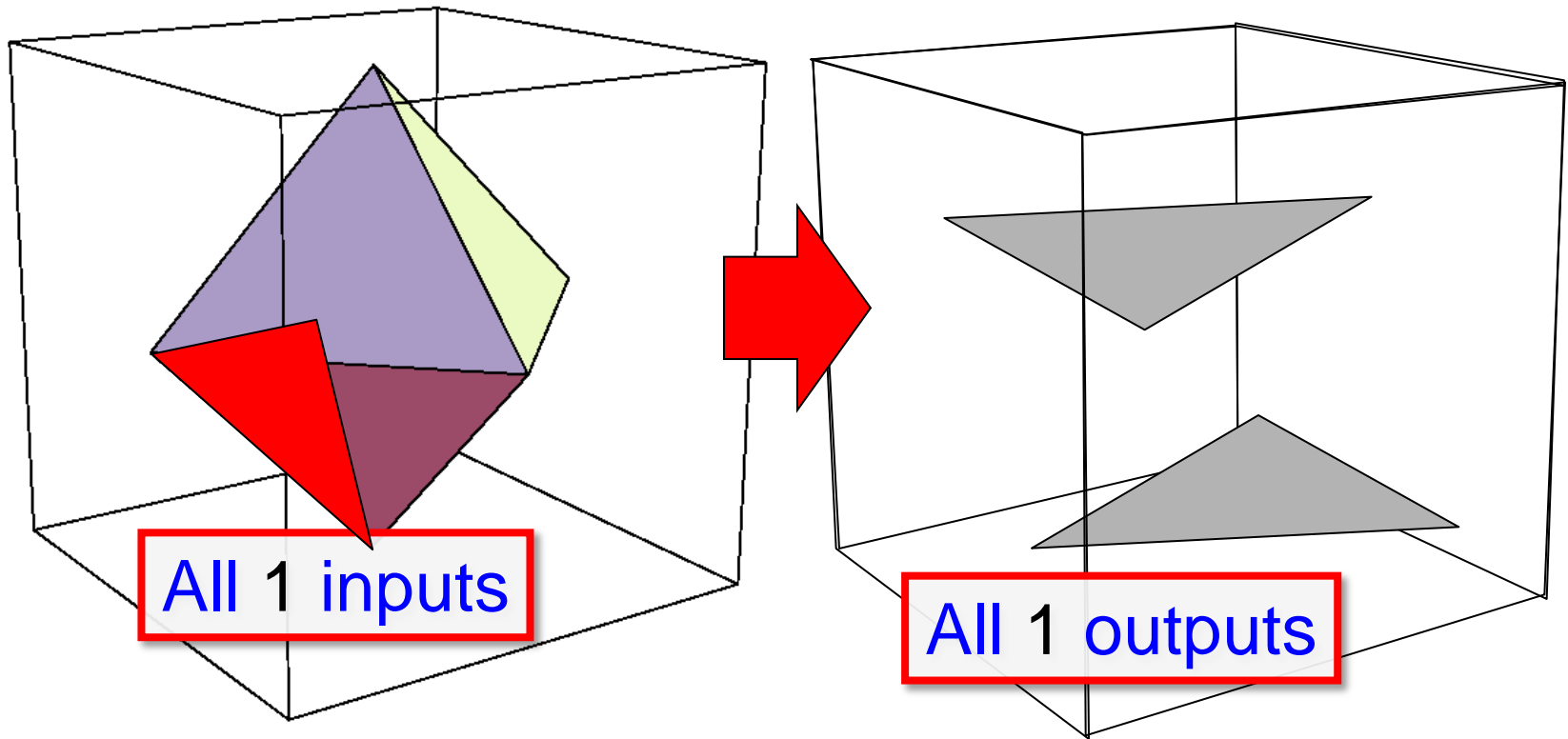
Output values all 0 or all 1

Two disconnected simplexes

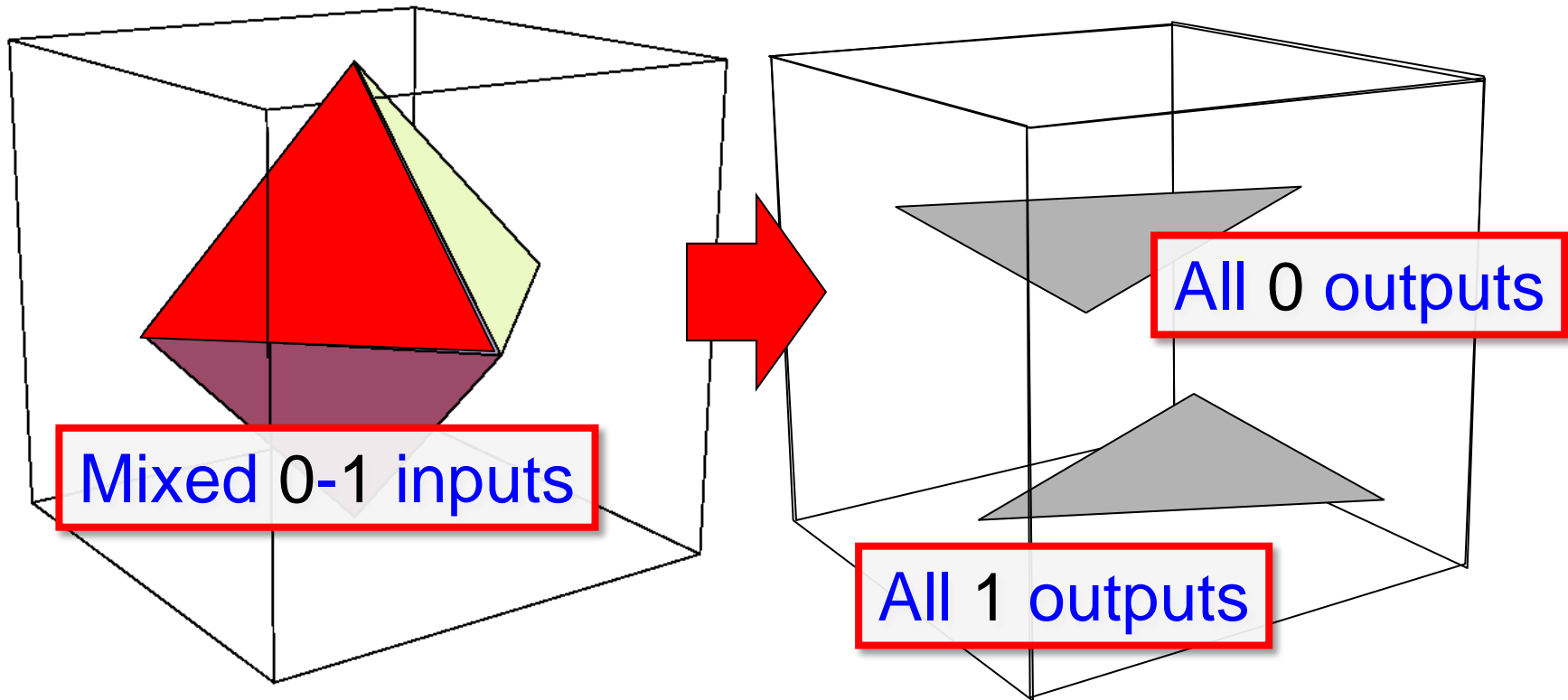
# Carrier Map for Consensus



# Carrier Map for Consensus

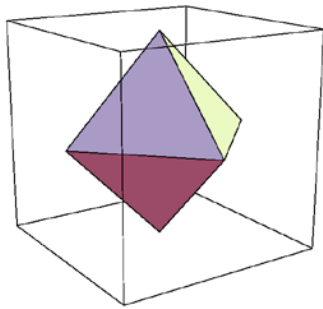


# Carrier Map for Consensus

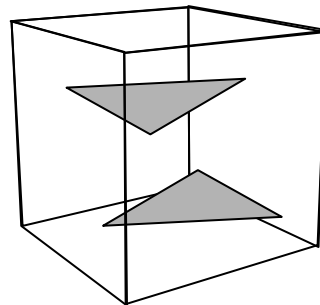


# Task Specification

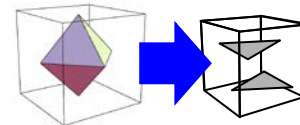
$(\mathcal{I}, \mathcal{O}, \Delta)$



Input complex



Output complex

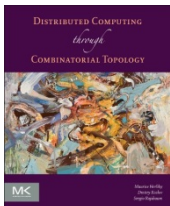
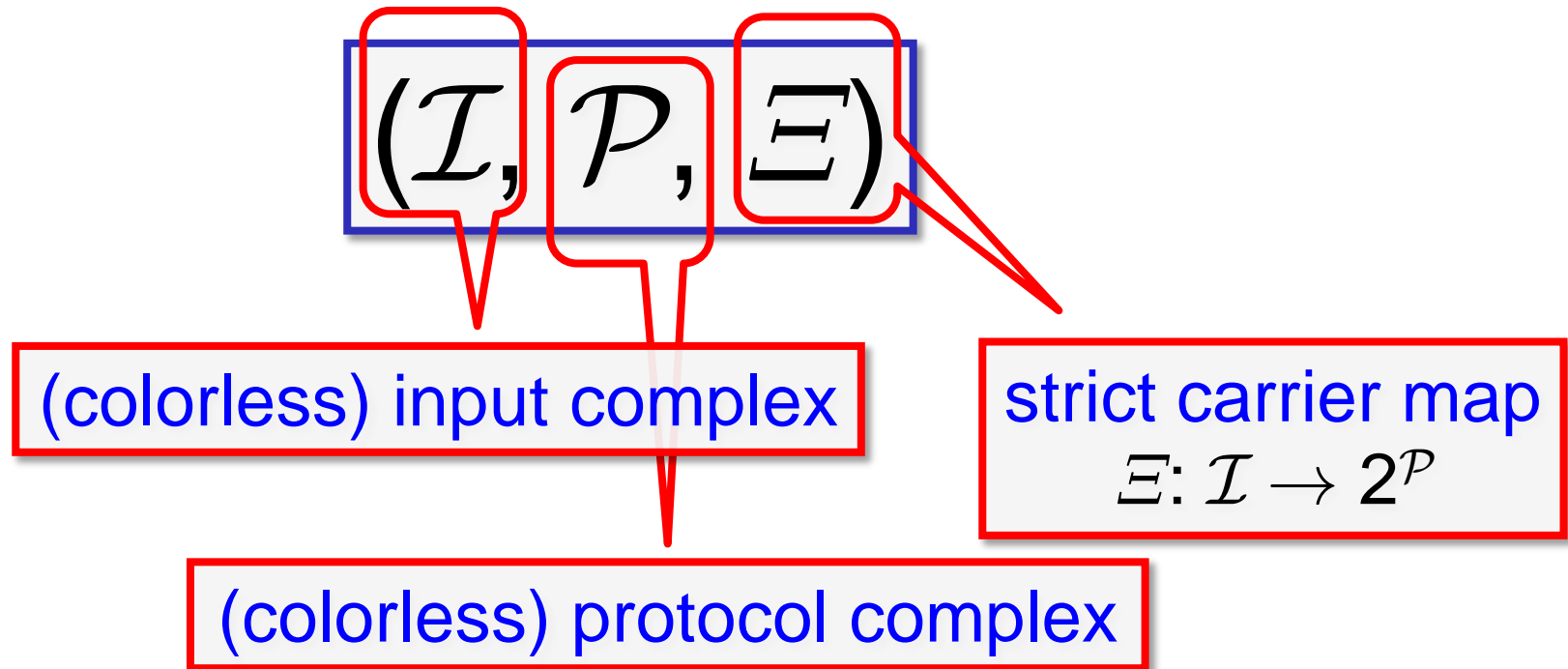


Carrier map

$$\Delta: \mathcal{I} \rightarrow 2^{\mathcal{O}}$$



# Colorless Tasks



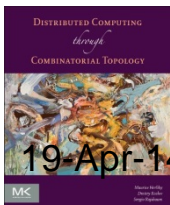
# Protocol Complex

Vertex: process name, view

all values read and written

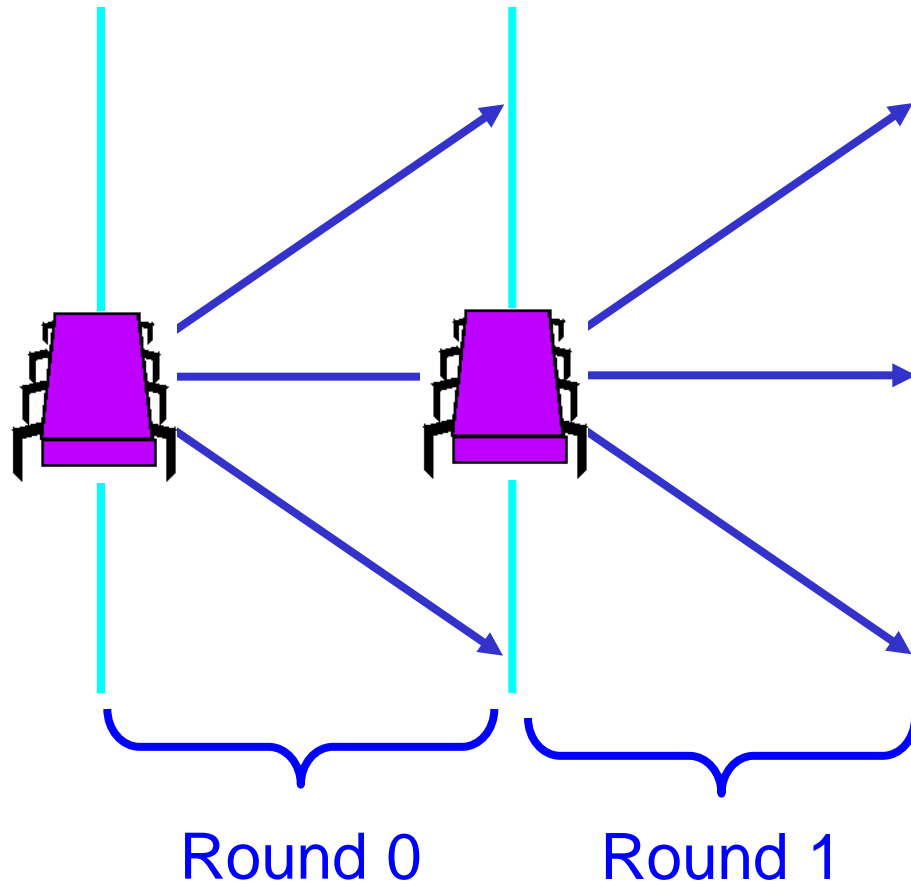
Simplex: compatible set of views

Each execution defines a simplex

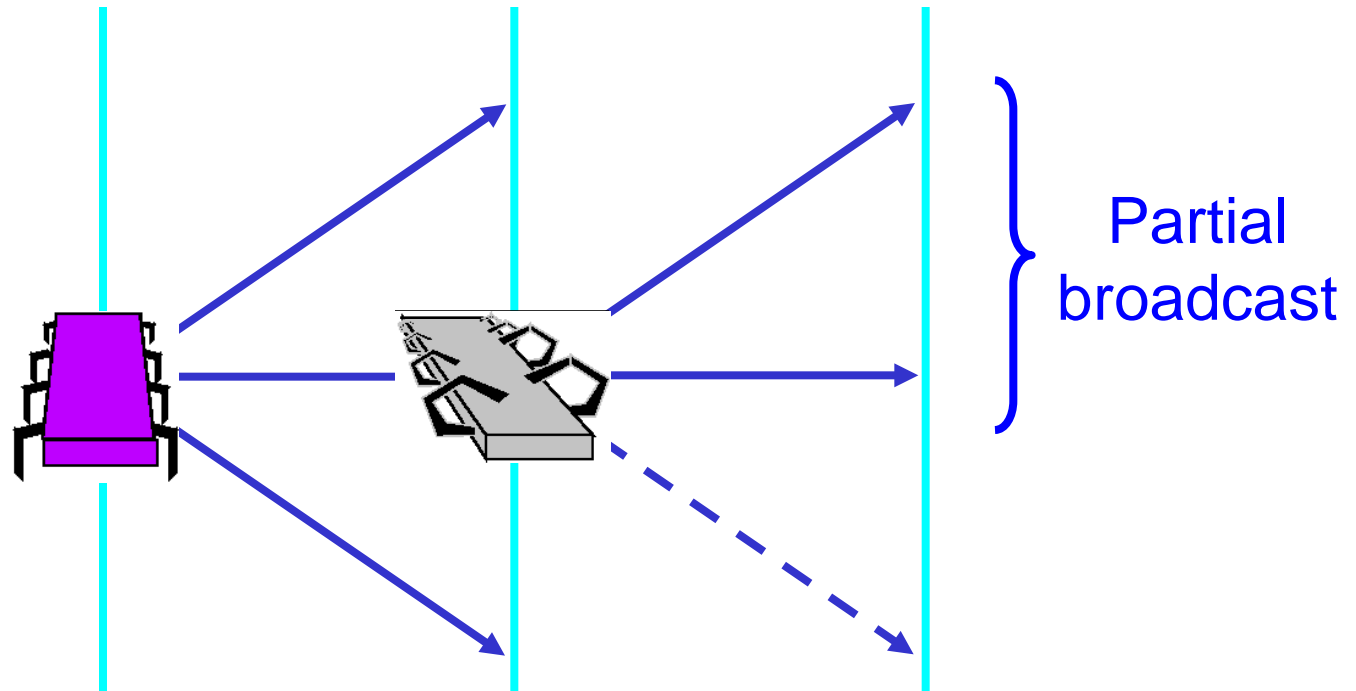


19-Apr-14

# Example: Synchronous Message-Passing

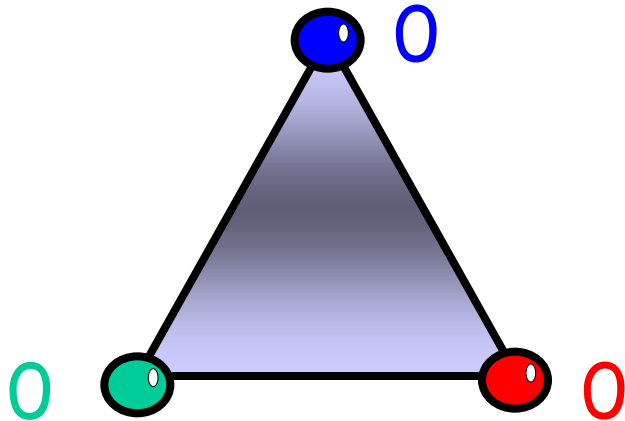


# Failures: Fail-Stop



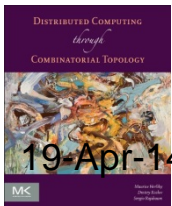
# Single Input: Round Zero

No messages sent

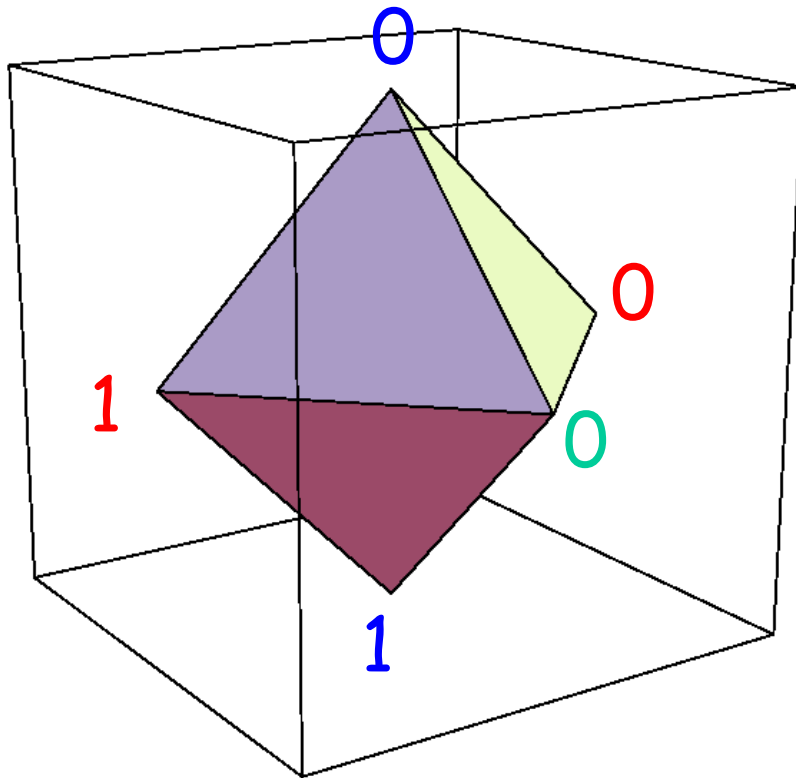


View is input value

Same as input simplex



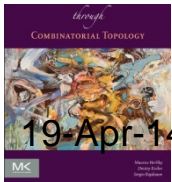
# Round Zero Protocol Complex



No messages sent

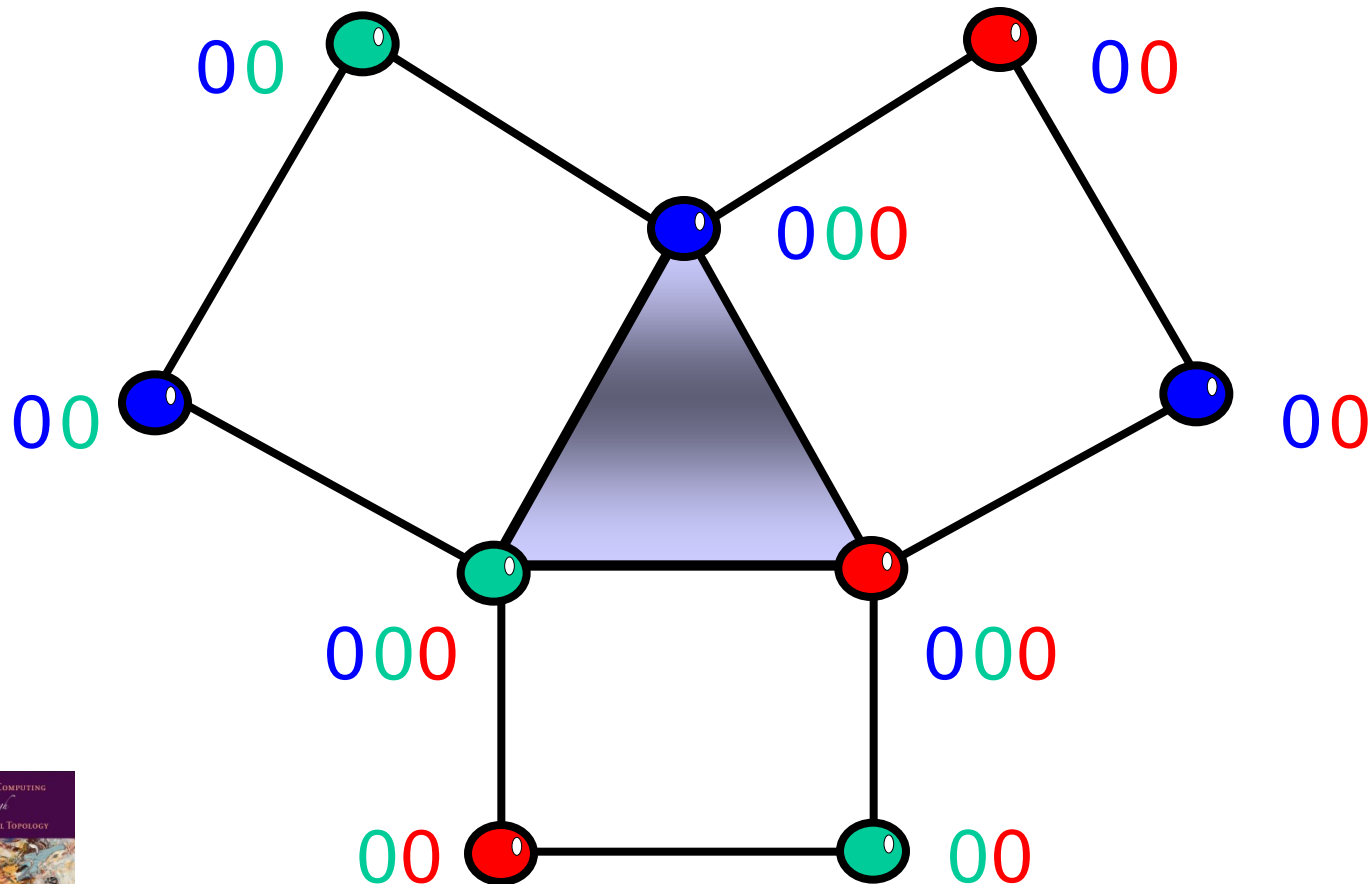
View is input value

Same as input complex

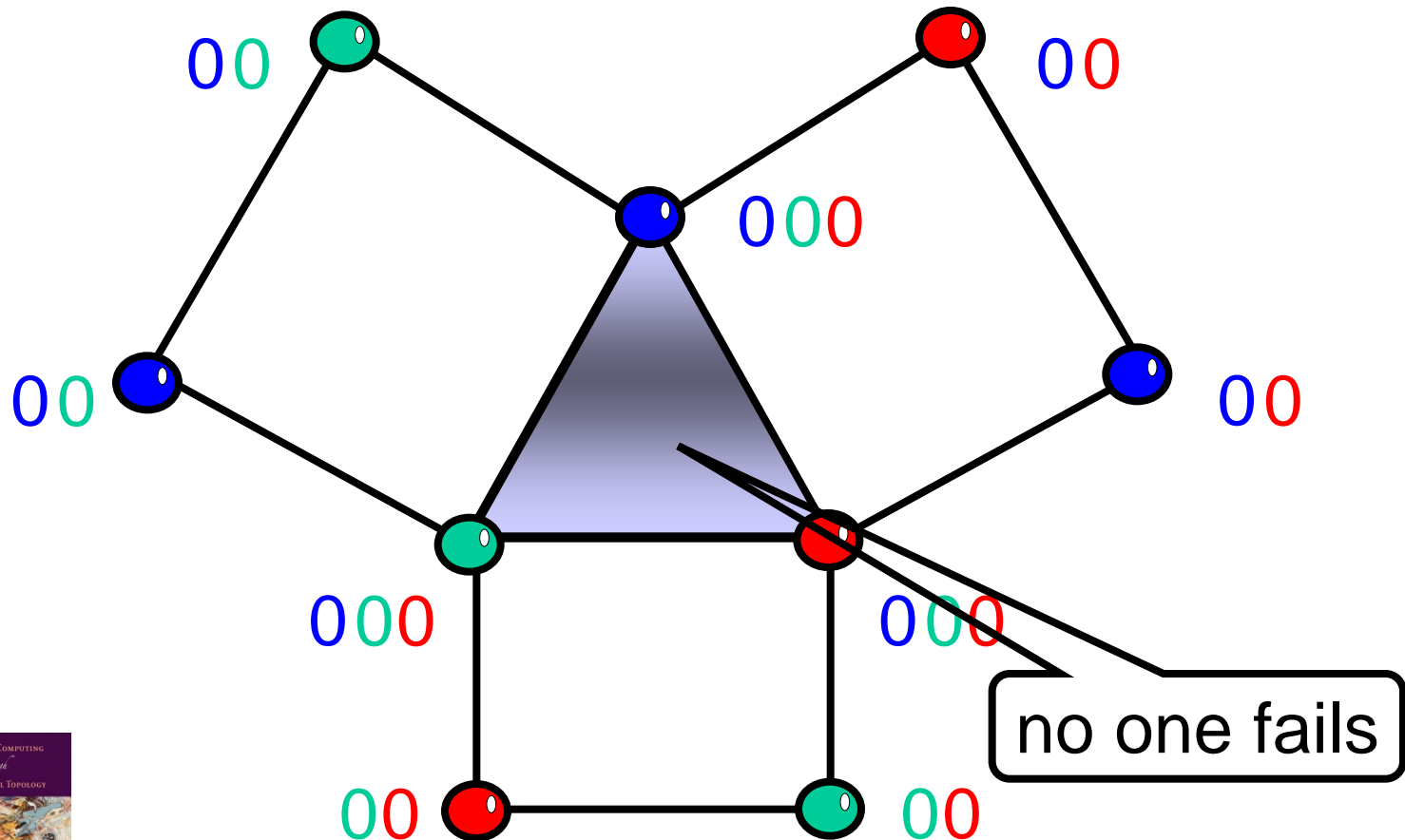


19-Apr-14

# Single Input: Round One



# Single Input: Round One



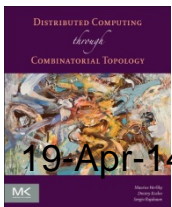
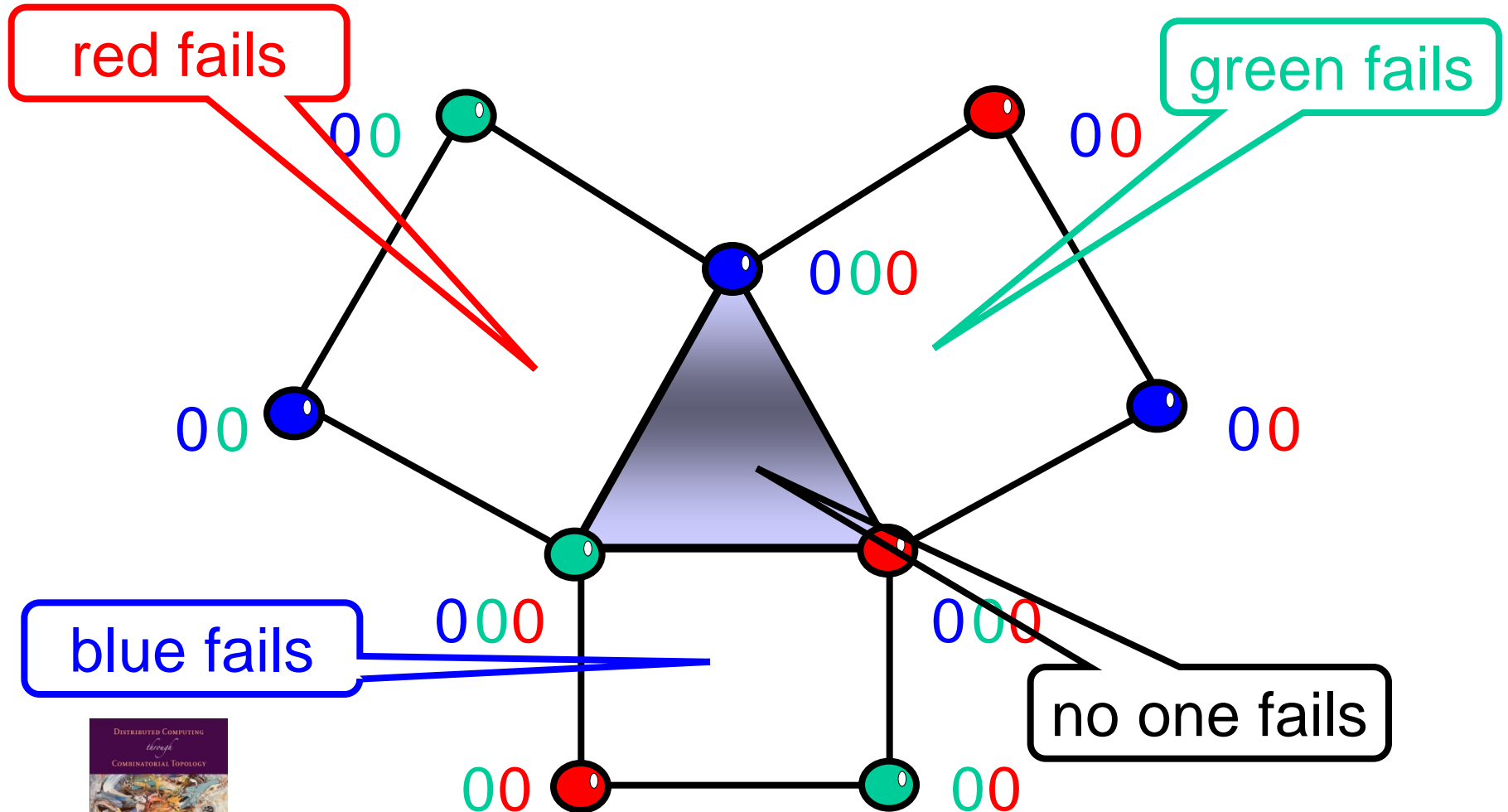


The diagram shows a network of nodes and edges. A central triangle is shaded blue. Nodes are colored red, green, or blue. Labels '00' are placed near nodes. A blue callout box points to a node with the text 'e fails'. A white callout box points to a node with the text 'no one fa'.

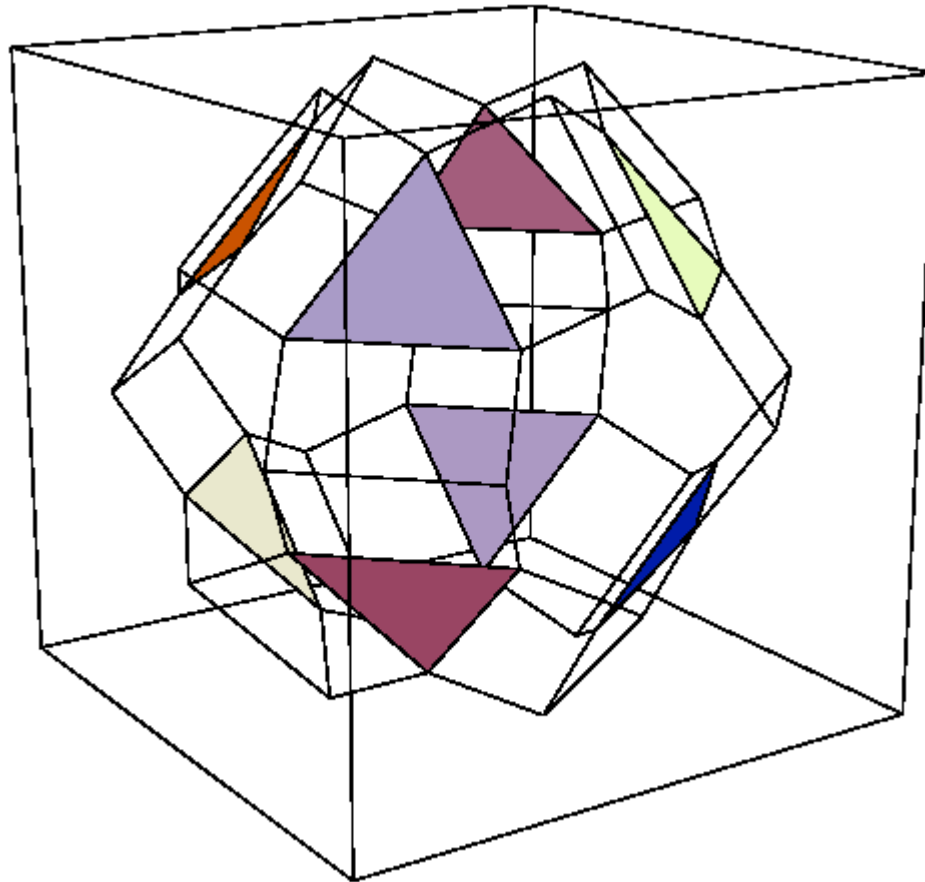
no one fails



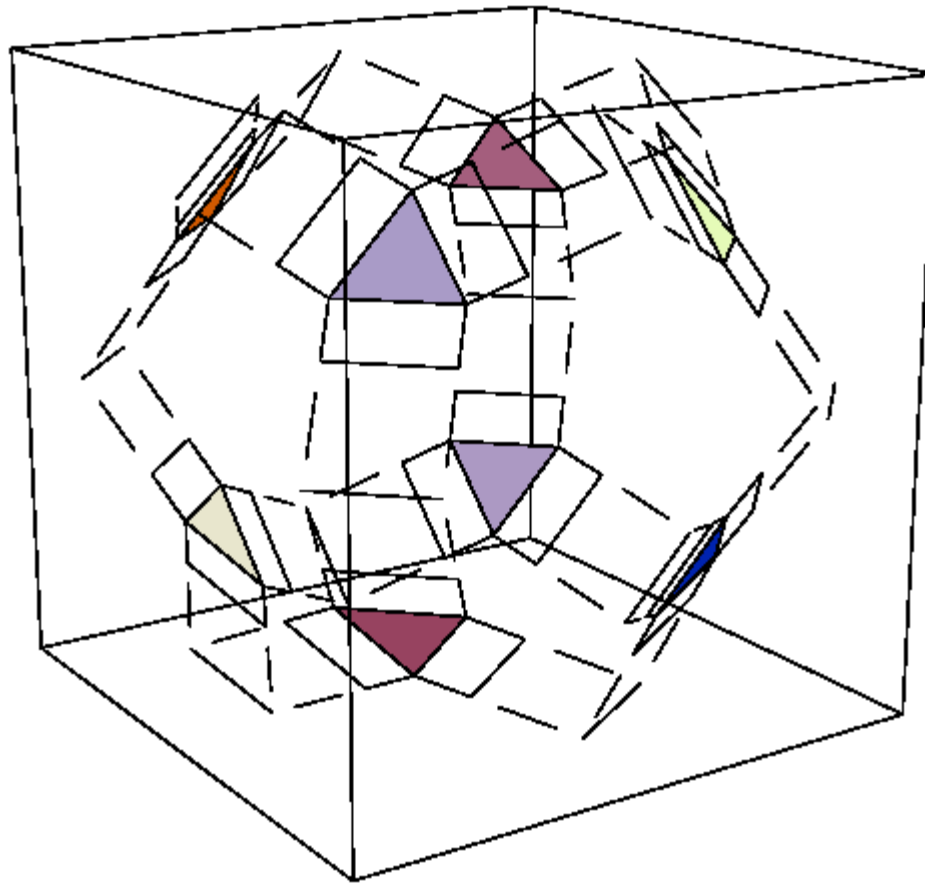
# Single Input: Round One



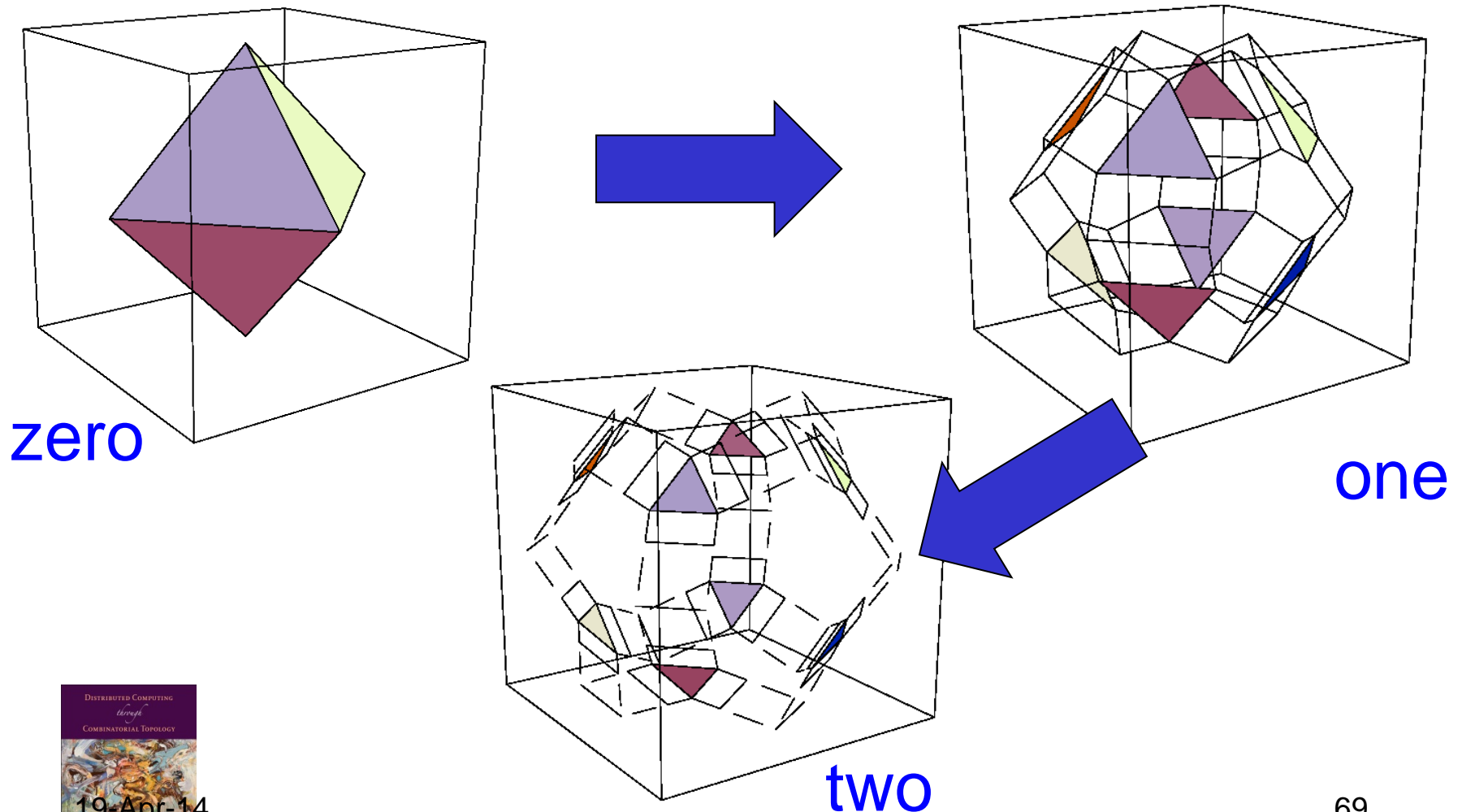
# Protocol Complex: Round One



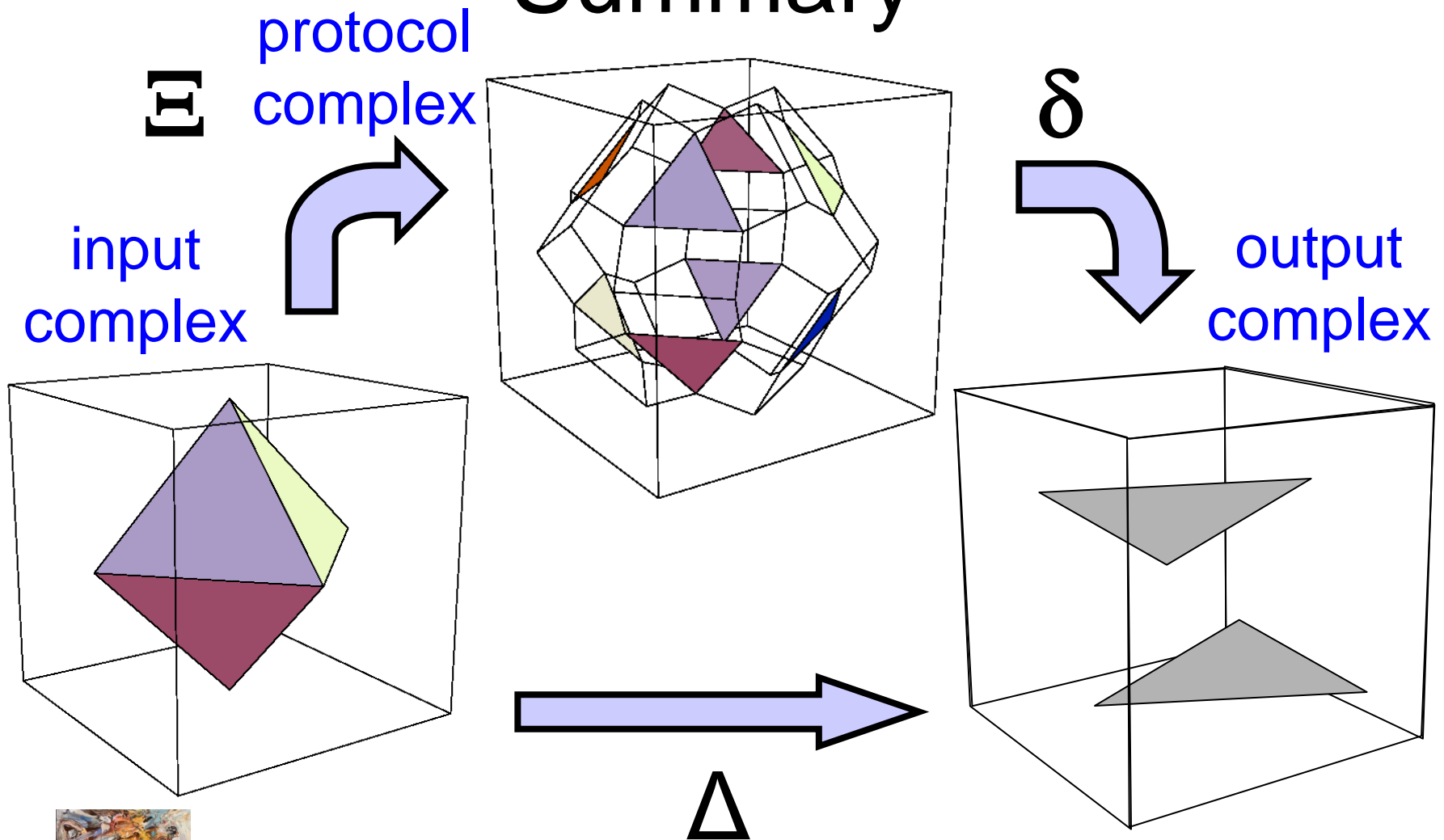
# Protocol Complex: Round Two



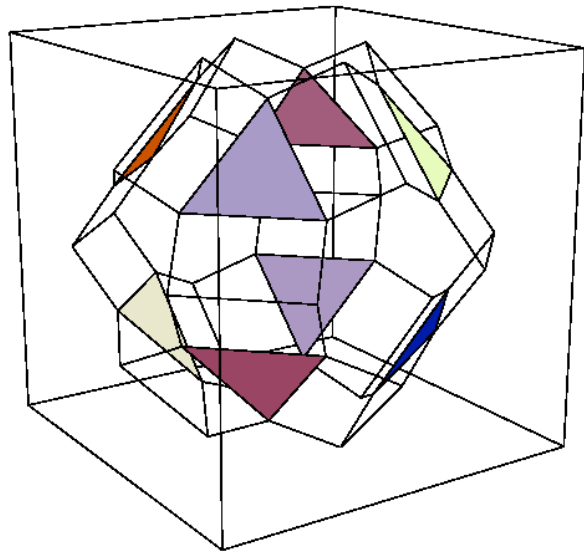
# Protocol Complex Evolution



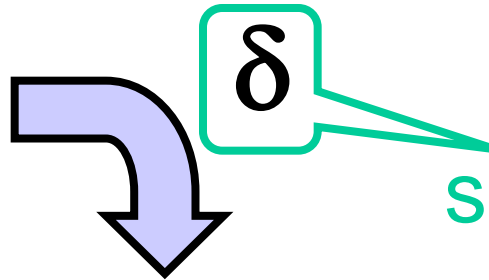
# Summary



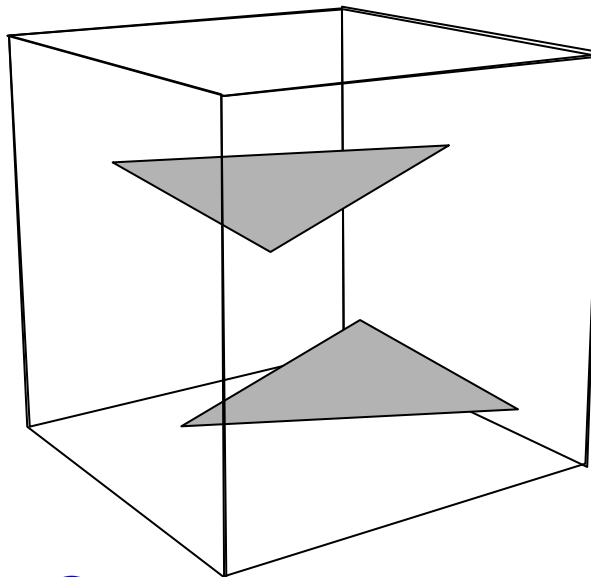
# Decision Map



Protocol complex

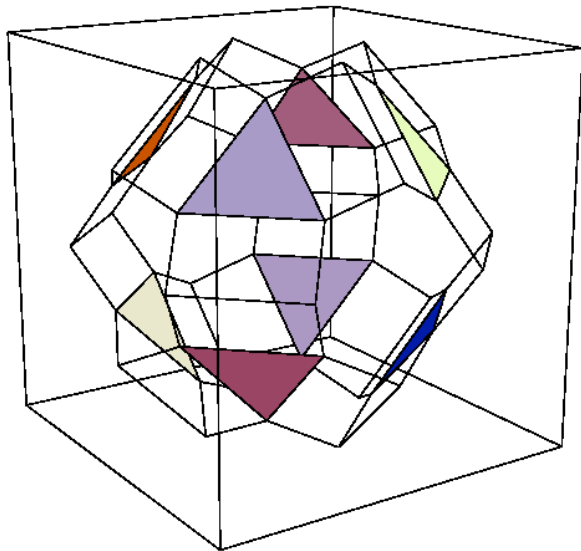


Simplicial map,  
sending simplexes  
to simplexes

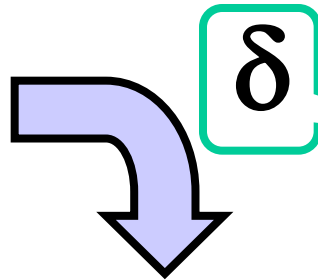


Output complex

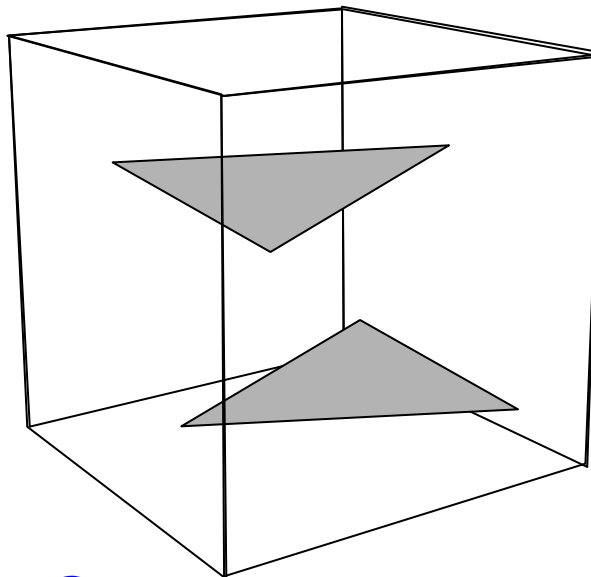
# Lower Bound Strategy



Protocol complex



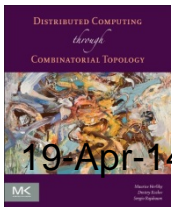
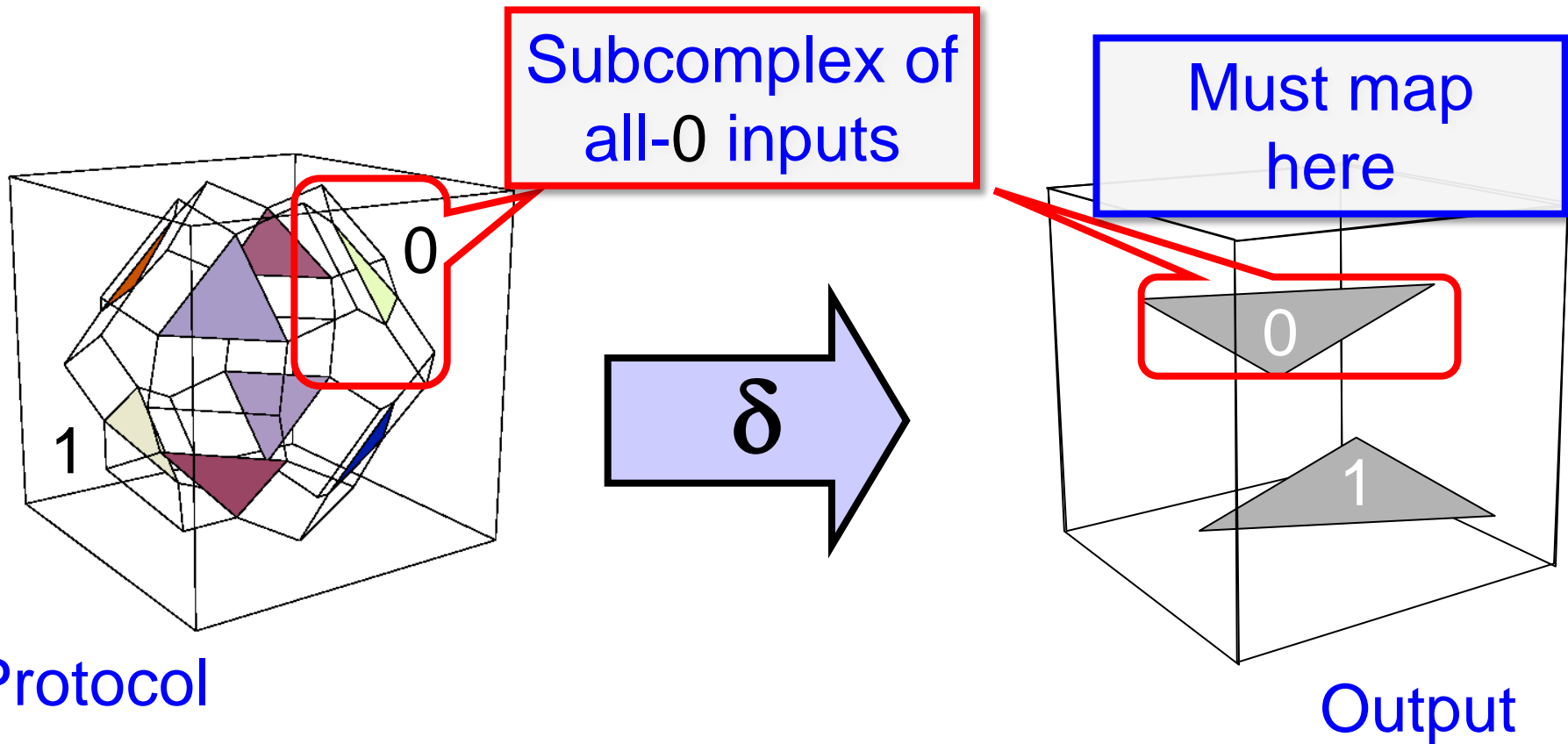
Find topological  
“obstruction” to  
this simplicial map



Output complex

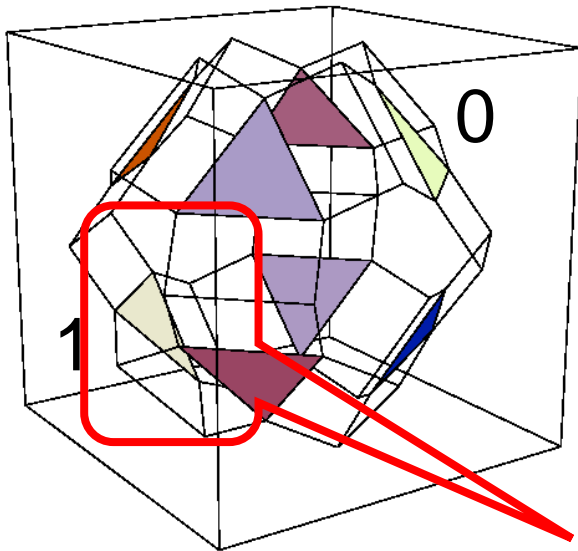


# Consensus Example

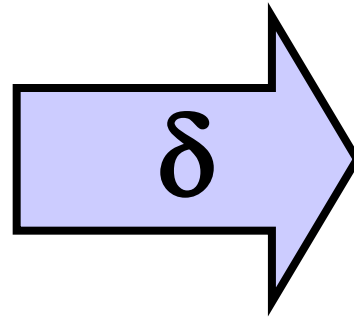


19-Apr-14

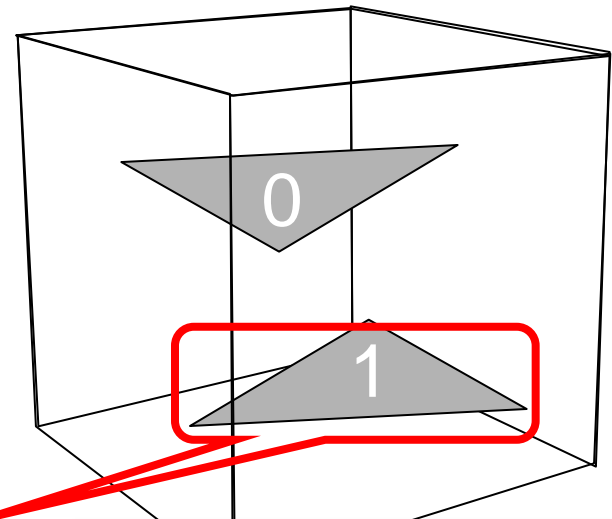
# Consensus Example



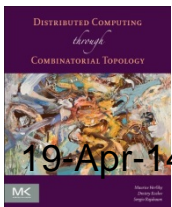
Protocol



Subcomplex of  
all-1 inputs



Must map  
here



19-Apr-14

# Consensus Example

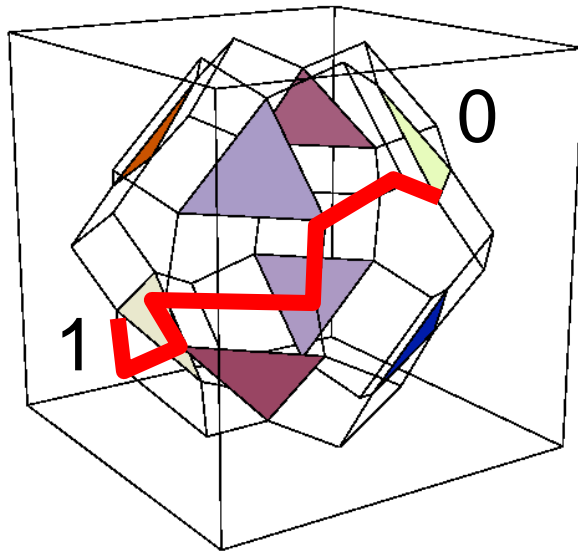
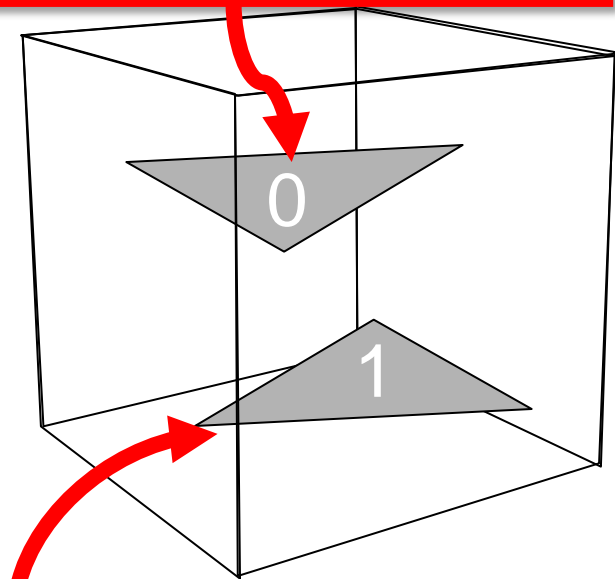
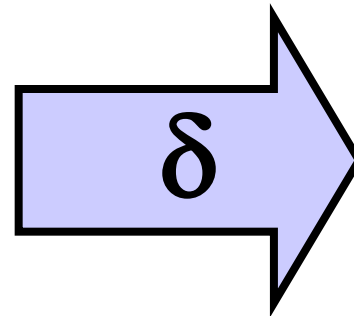


Image under  $\delta$  must start here ..

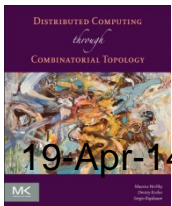
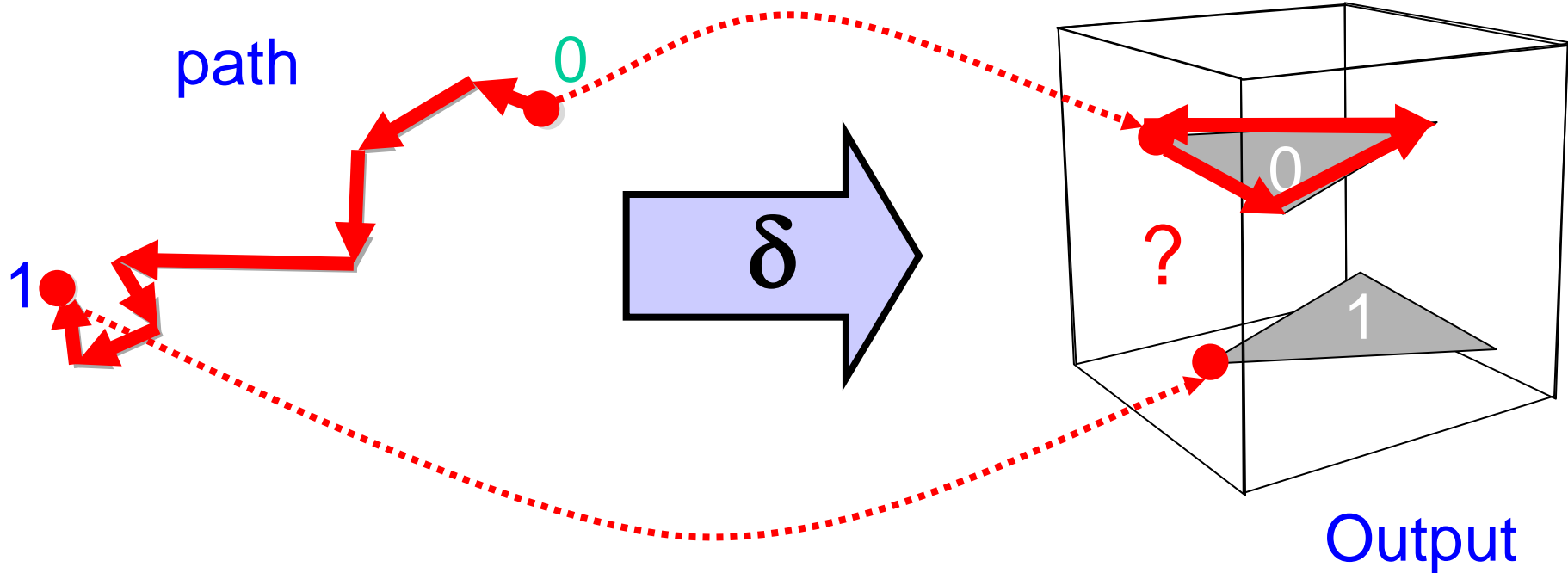


and end here

Output

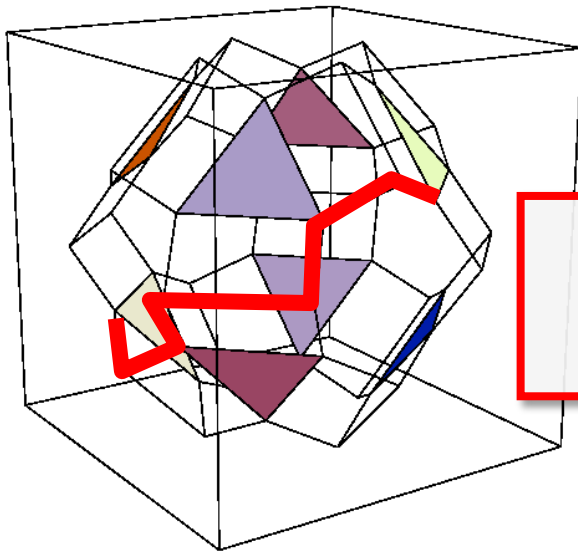
Path from  
“all-0” to “all-1”

# Consensus Example

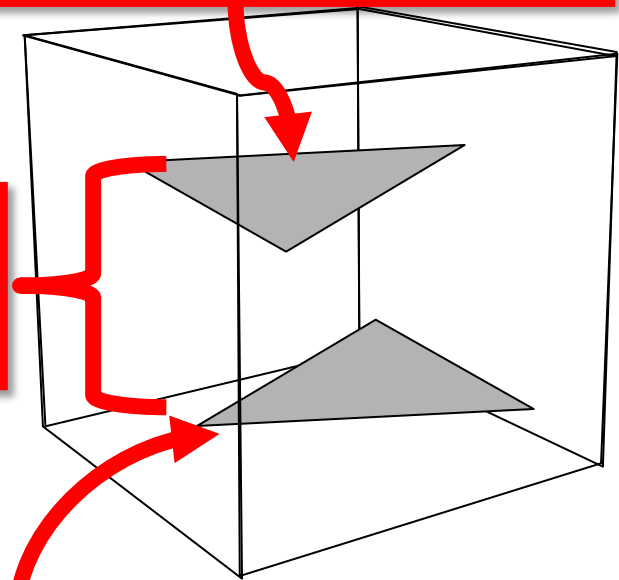


# Consensus Example

Image under  $\delta$  must start here ..



But this “hole” is an obstruction

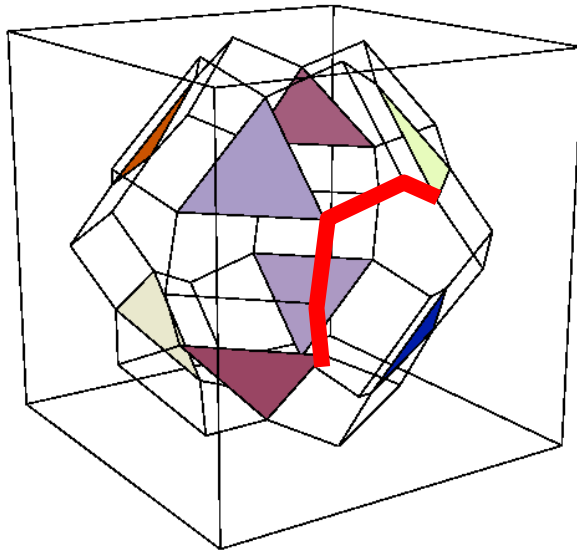


Path from  
“all-0” to “all-1”

and end here

Output

# Conjecture



A protocol cannot  
solve consensus  
if its complex is  
***path-connected***

**Model-independent!**

# If Adversary keeps Protocol Complex path-connected ...

Forever ...

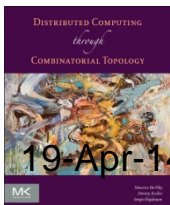
Consensus is *impossible*

For  $r$  rounds ...

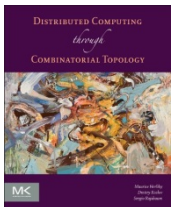
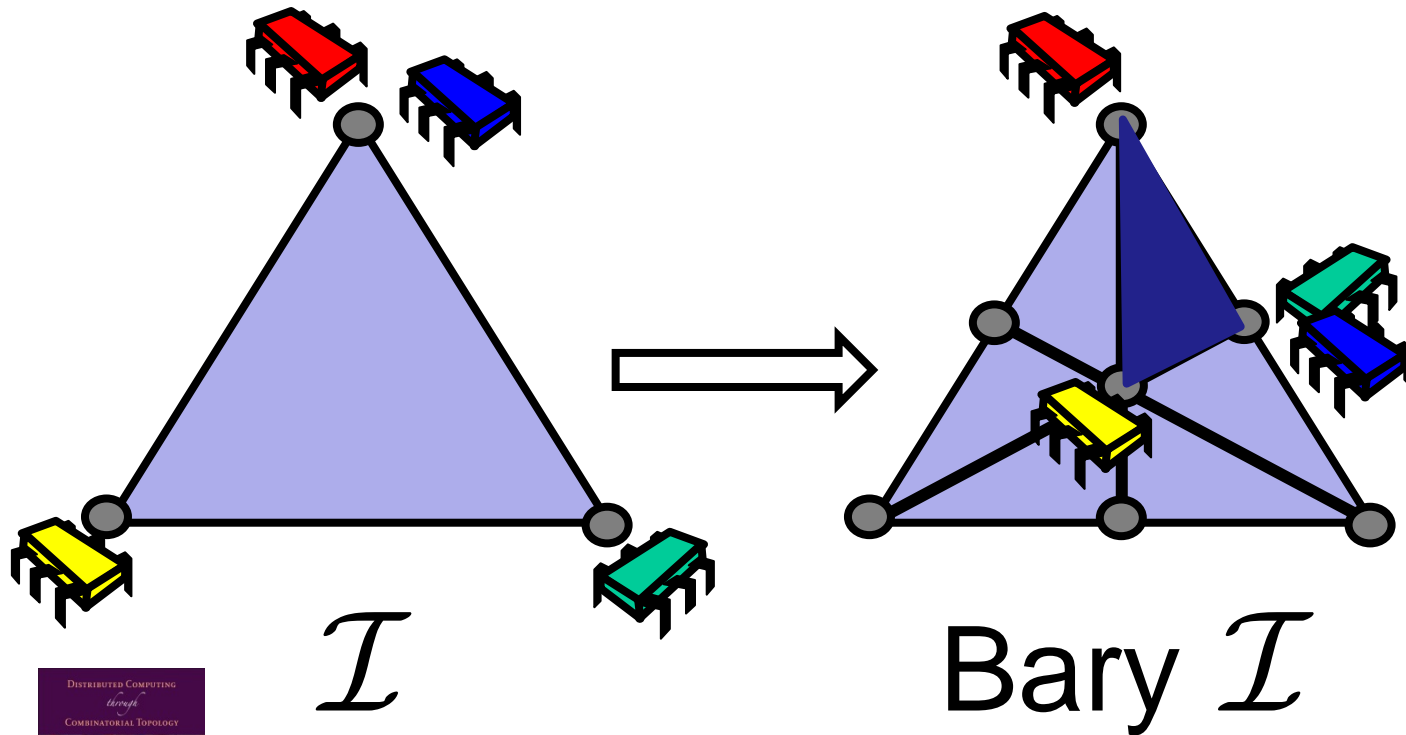
A *round-complexity* lower bound

For time  $t$  ...

A *time-complexity* lower bound

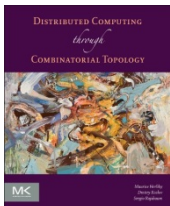
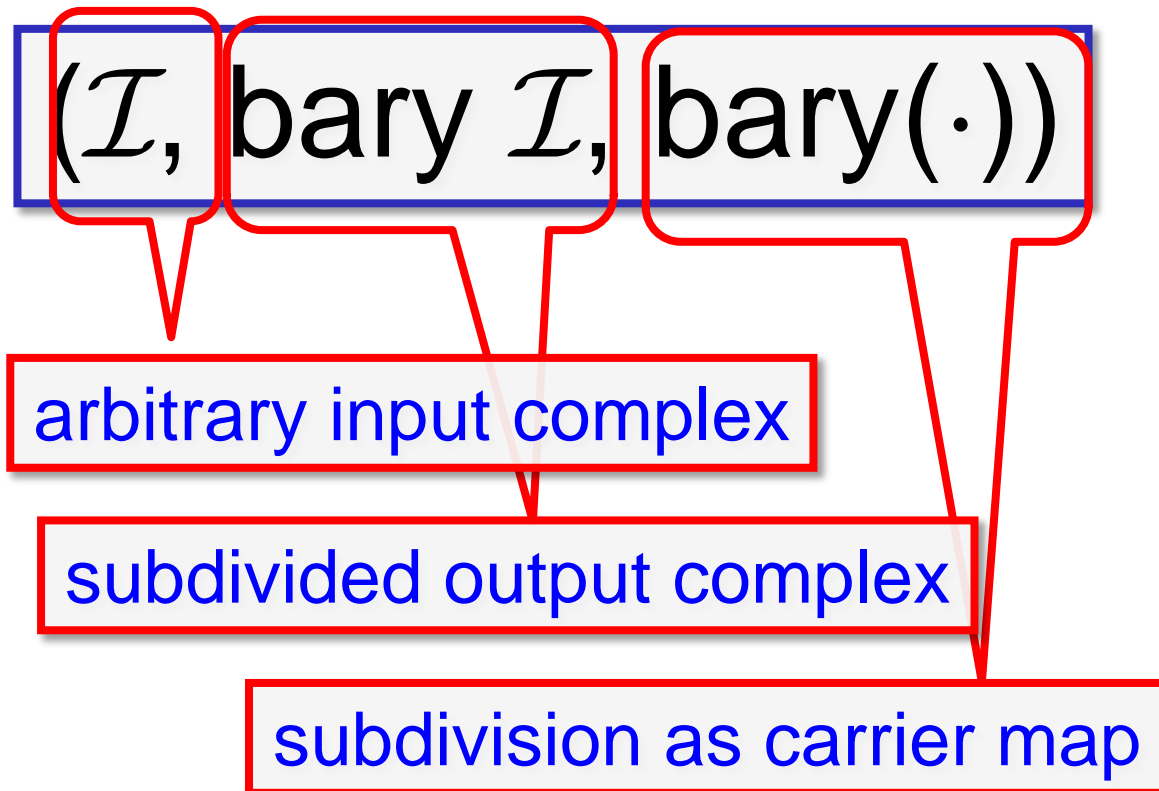


# Barycentric Agreement





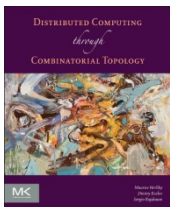
# Barycentric Agreement



If There are  $n$  Processes

$(\mathcal{I}, \text{bary} \text{skel}^n \mathcal{I}, \text{bary}^\circ \text{skel}^n)$

Inputs only from  $n$ -skeleton of input complex



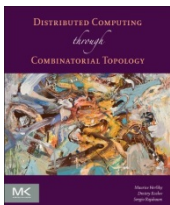
# Theorem

A one-layer immediate snapshot protocol solves the  $n$ -process barycentric agreement task  $(\mathcal{I}, \text{bary skel}^n \mathcal{I}, \text{bary}^\circ \text{skel}^n)$

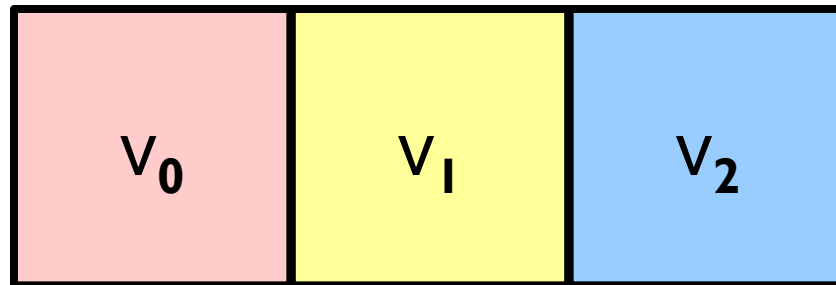
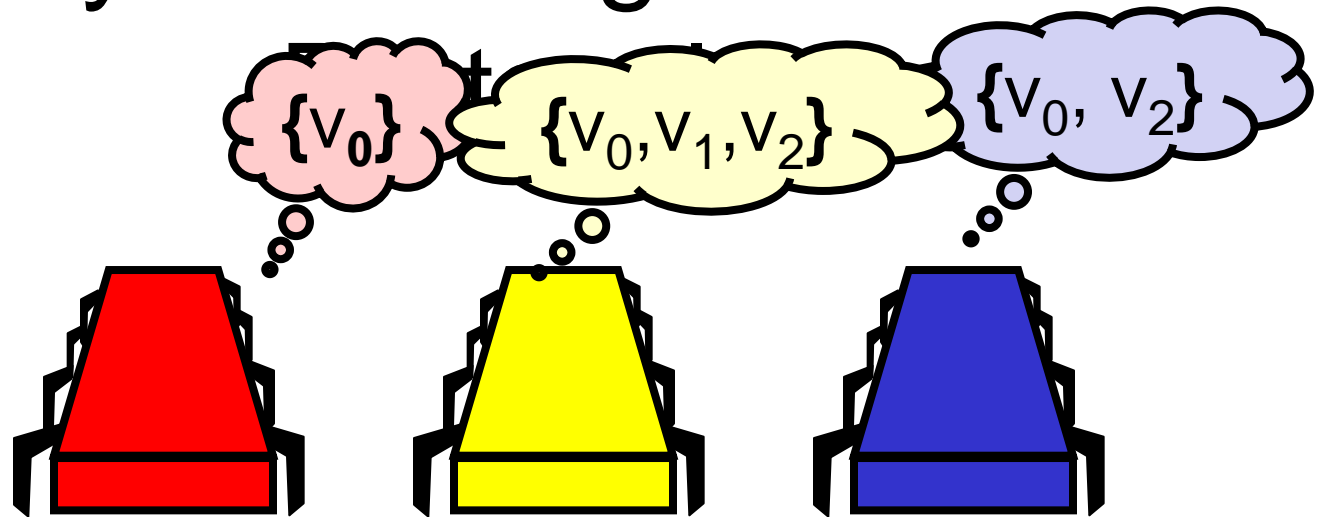
Proof

All input simplices belong to  $\text{skel}^n \mathcal{I}$

Immediate snapshot results are ordered

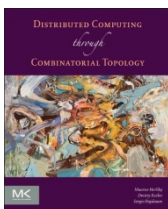
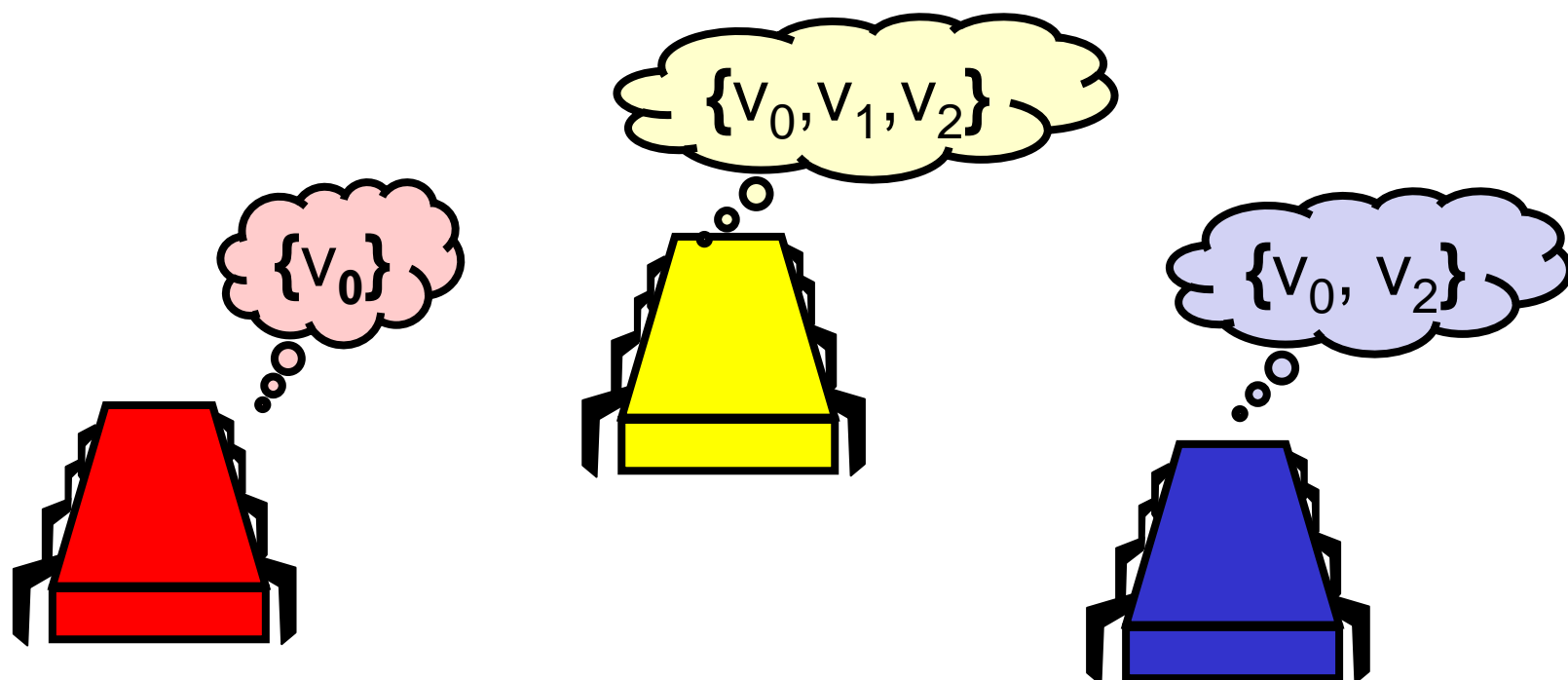


# Barycentric Agreement



Snapshots are ordered

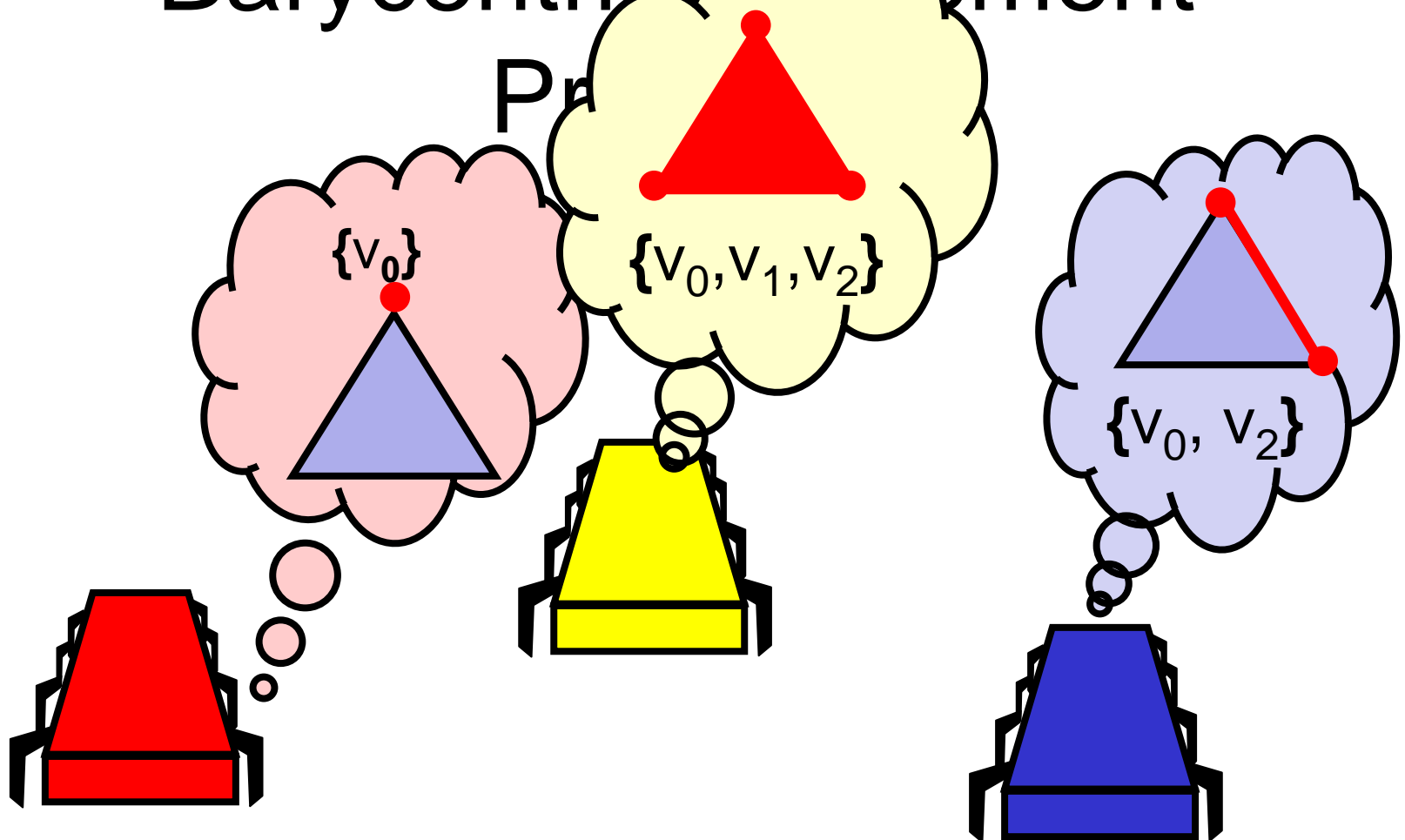
# Barycentric Agreement Protocol



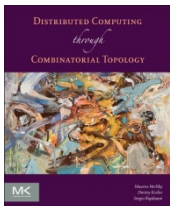
Each view is a face of  $\sigma$

# Barycentric Arrangement

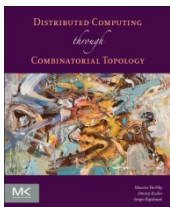
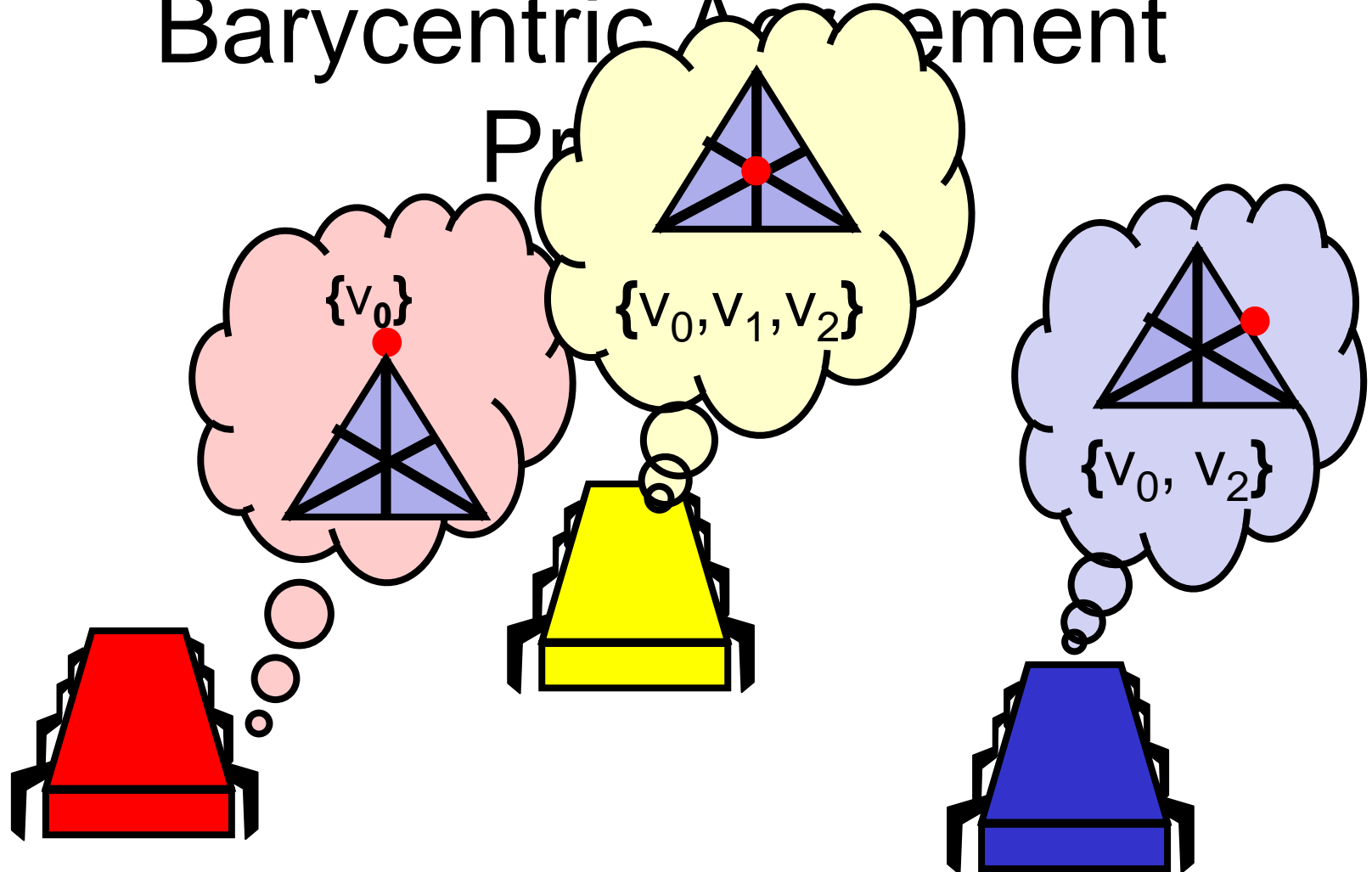
Pr



Each view is a face of  $\sigma$



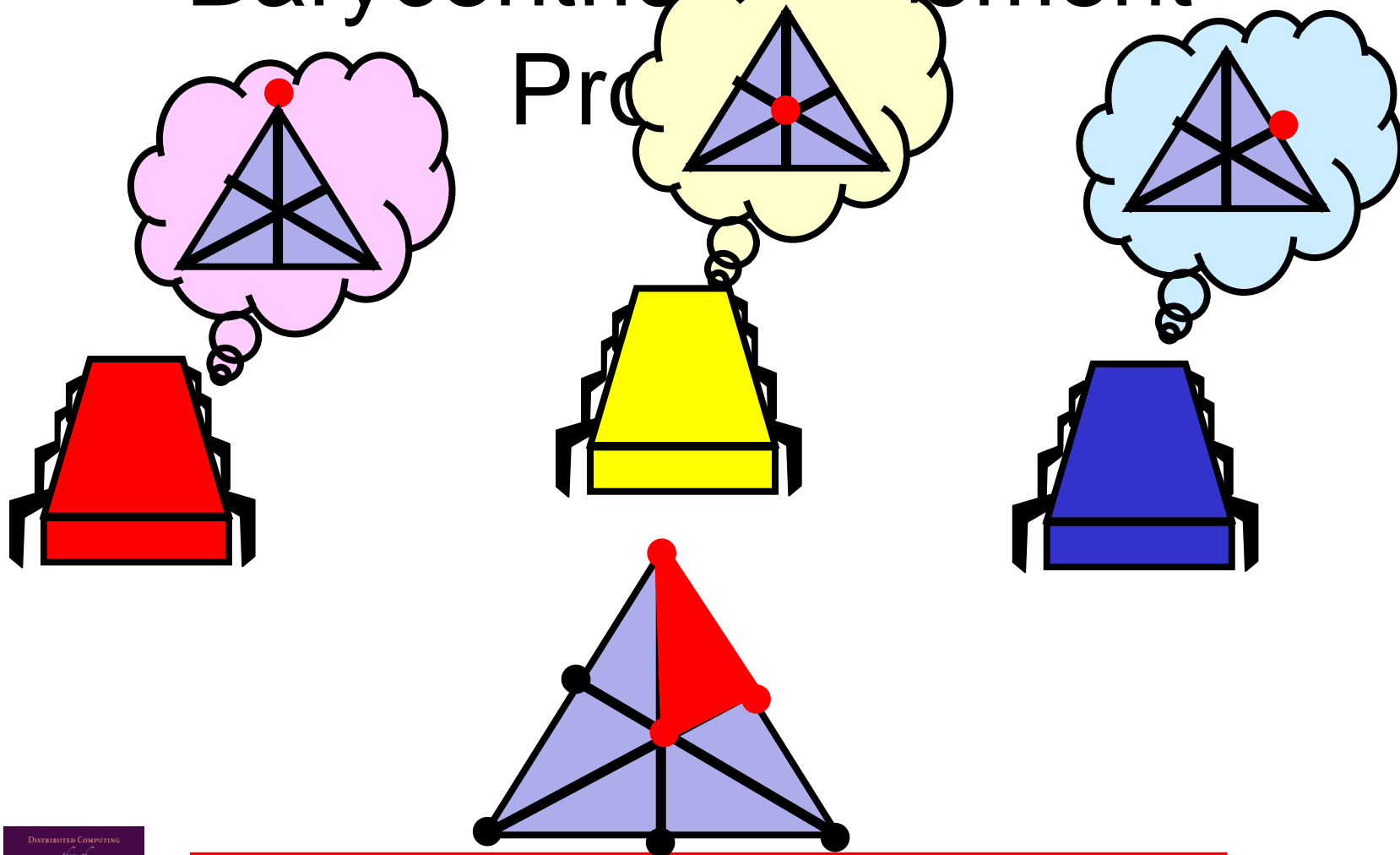
# Barycentric Arrangement



Each face of  $\sigma$  is a vertex of Bary  $\sigma$

# Barycentric Arrangement

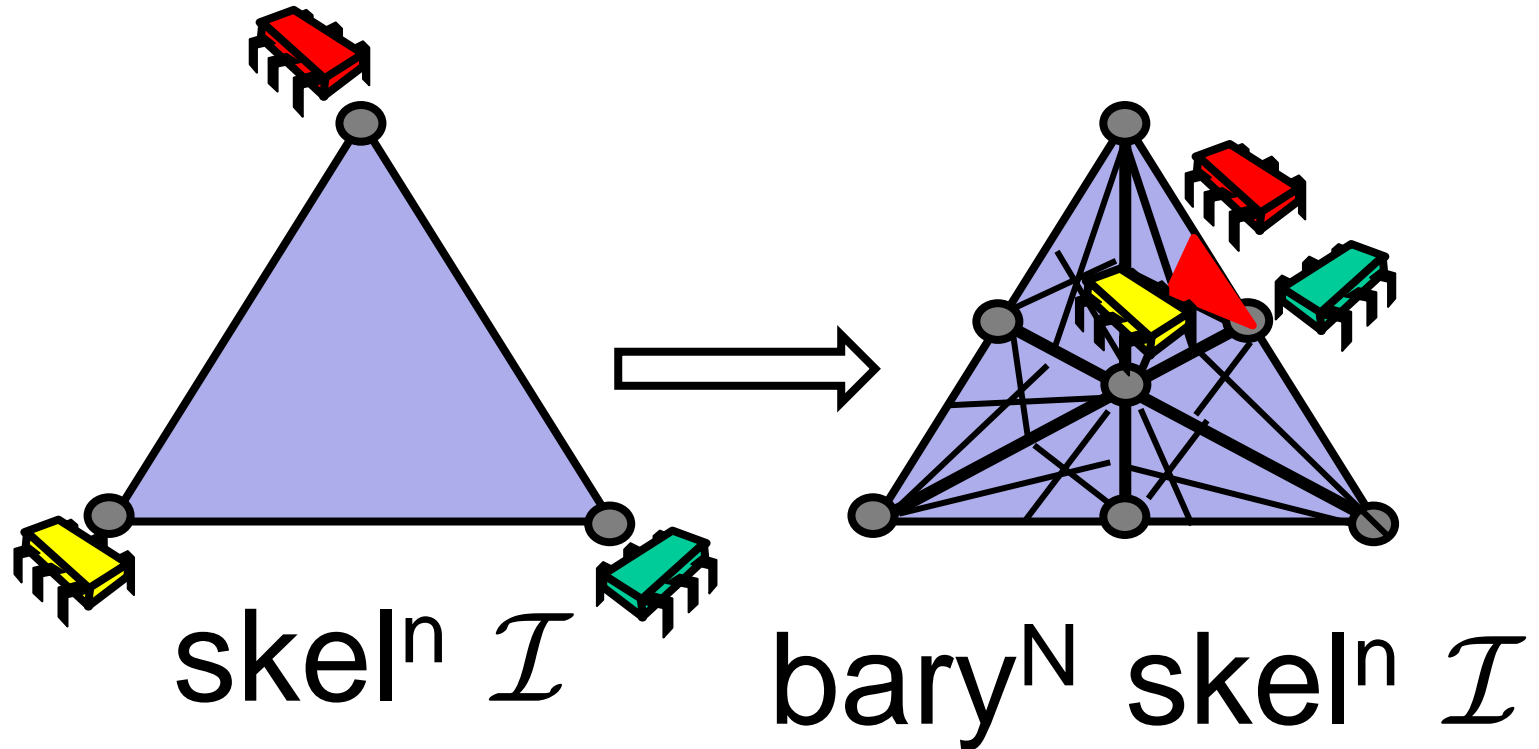
## Process



Ordered faces  $\Rightarrow$  simplex of Bary  $\sigma$

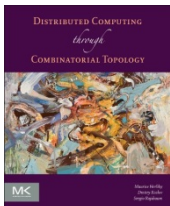


# Iterated Barycentric Agreement

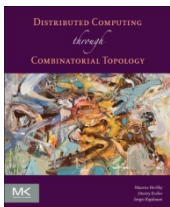
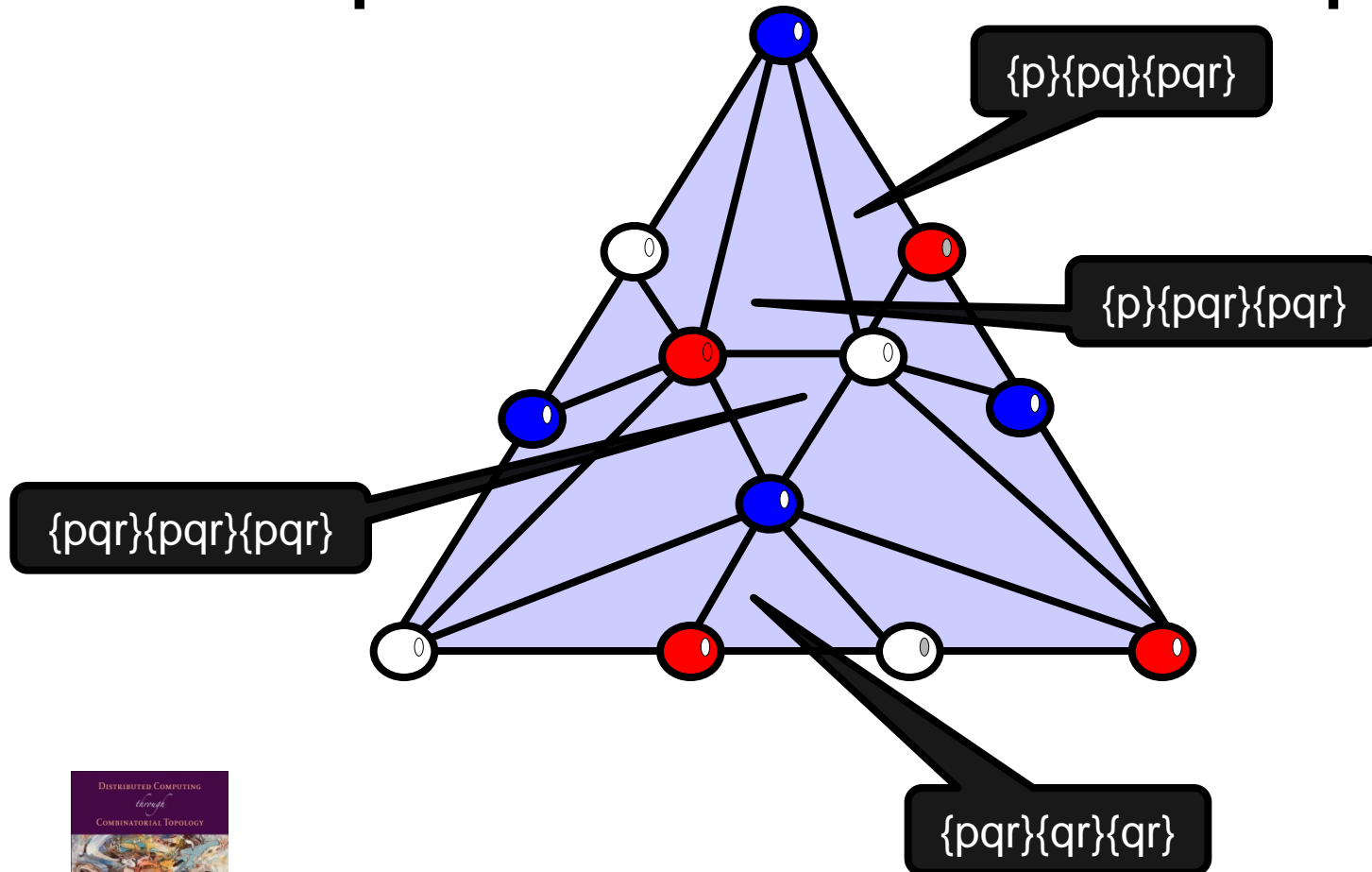


# Iterated Barycentric Agreement

$$(\mathcal{I}, \text{bary}^N \text{skel}^n \mathcal{I}, \text{bary}^{N^o} \text{skel}^n)$$

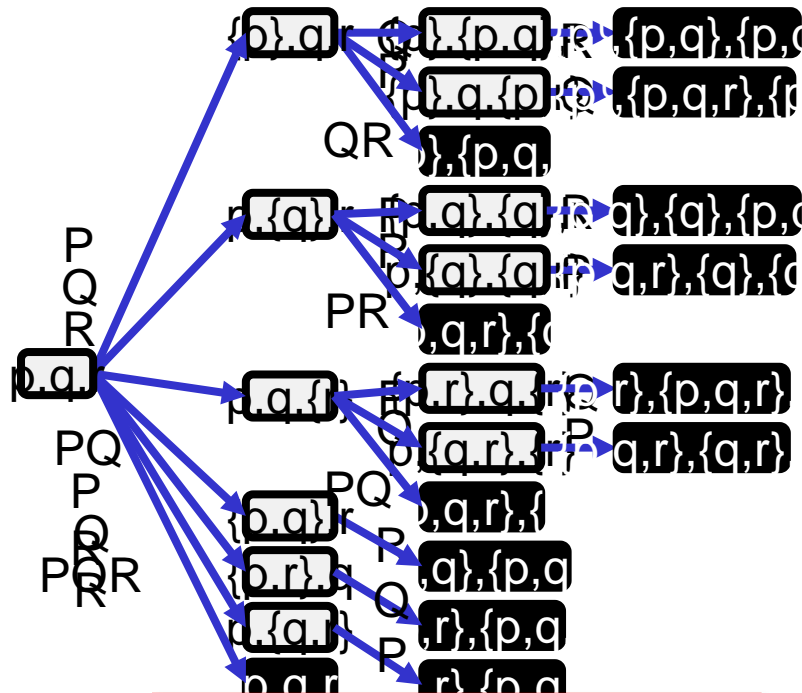


# One-Layer Immediate Snapshot Protocol Complex

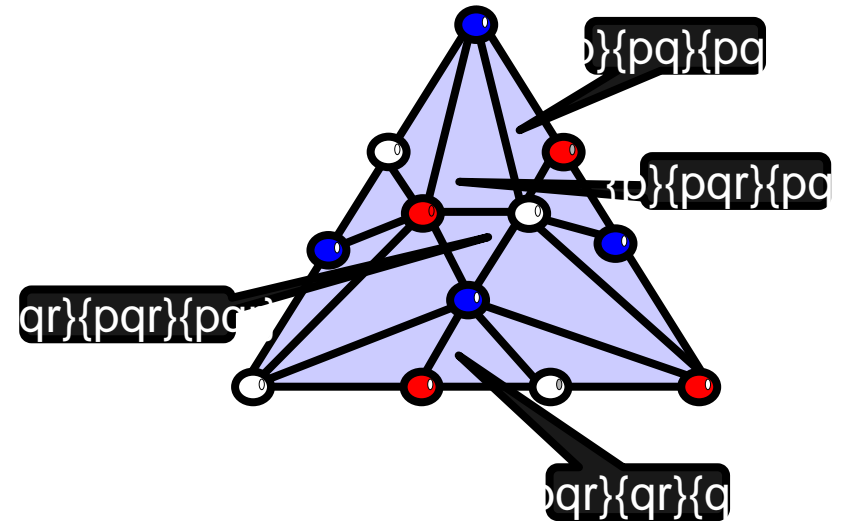


Distributed Computing through  
Combinatorial Topology

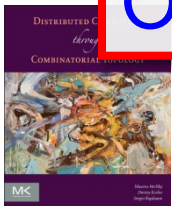
# Compare Views



Operational view



Combinatorial view



# Compositions

Given protocols

$$(\mathcal{I}, \mathcal{P}, \Xi)$$

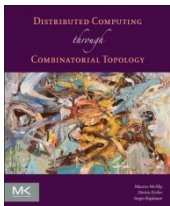
$$(\mathcal{I}, \mathcal{P}', \Xi')$$

where  $\mathcal{P} \subseteq \mathcal{I}$

their *composition* is

$$(\mathcal{I}, \mathcal{P}'', \Xi'')$$

where  $\Xi'' = \Xi' \circ \Xi$  and  $\mathcal{P}'' = \Xi''(\mathcal{I})$

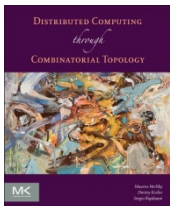


# Road Map

Operational Model

Combinatorial Model

Main Theorem



Distributed Computing through  
Combinatorial Topology

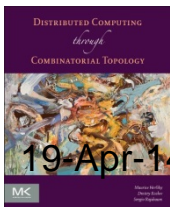
# Fundamental Theorem

## Theorem

$(\mathcal{I}, \mathcal{O}, \Delta)$  has a wait-free  $(n+1)$ -process layered protocol iff there is a continuous map

$$f: |\text{skel}^n \mathcal{I}| \rightarrow |\mathcal{O}| \dots$$

carried by  $\Delta$



19-Apr-14

Distributed Computing through  
Combinatorial Topology

# Proof Outline

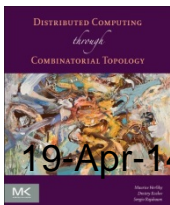
Lemma

If ...

there is a WF layered protocol for  
 $(\mathcal{I}, \mathcal{O}, \Delta) \dots$

then ...

there is a continuous  
 $f: |\text{skel}^n \mathcal{I}| \rightarrow |\mathcal{O}|$  carried by  $\Delta$ .

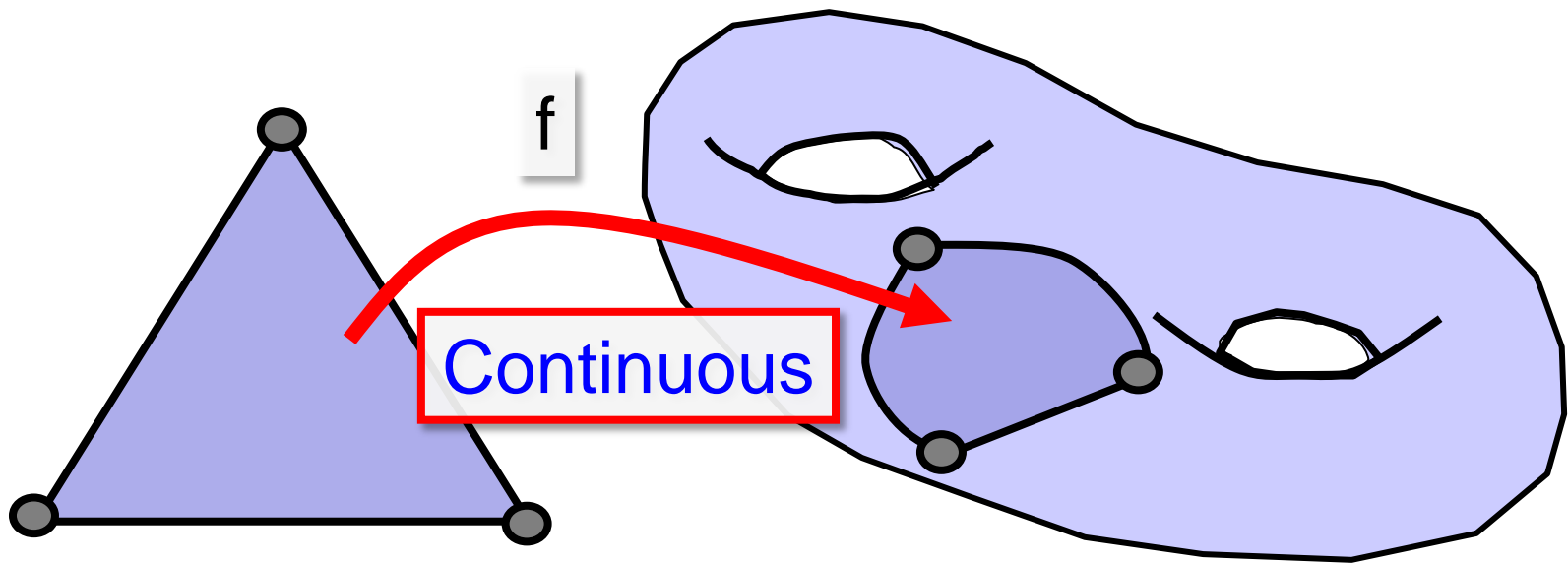


19-Apr-14



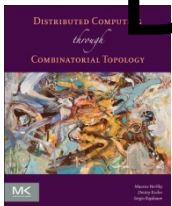
# Map $\Rightarrow$ Protocol

Hypothesis

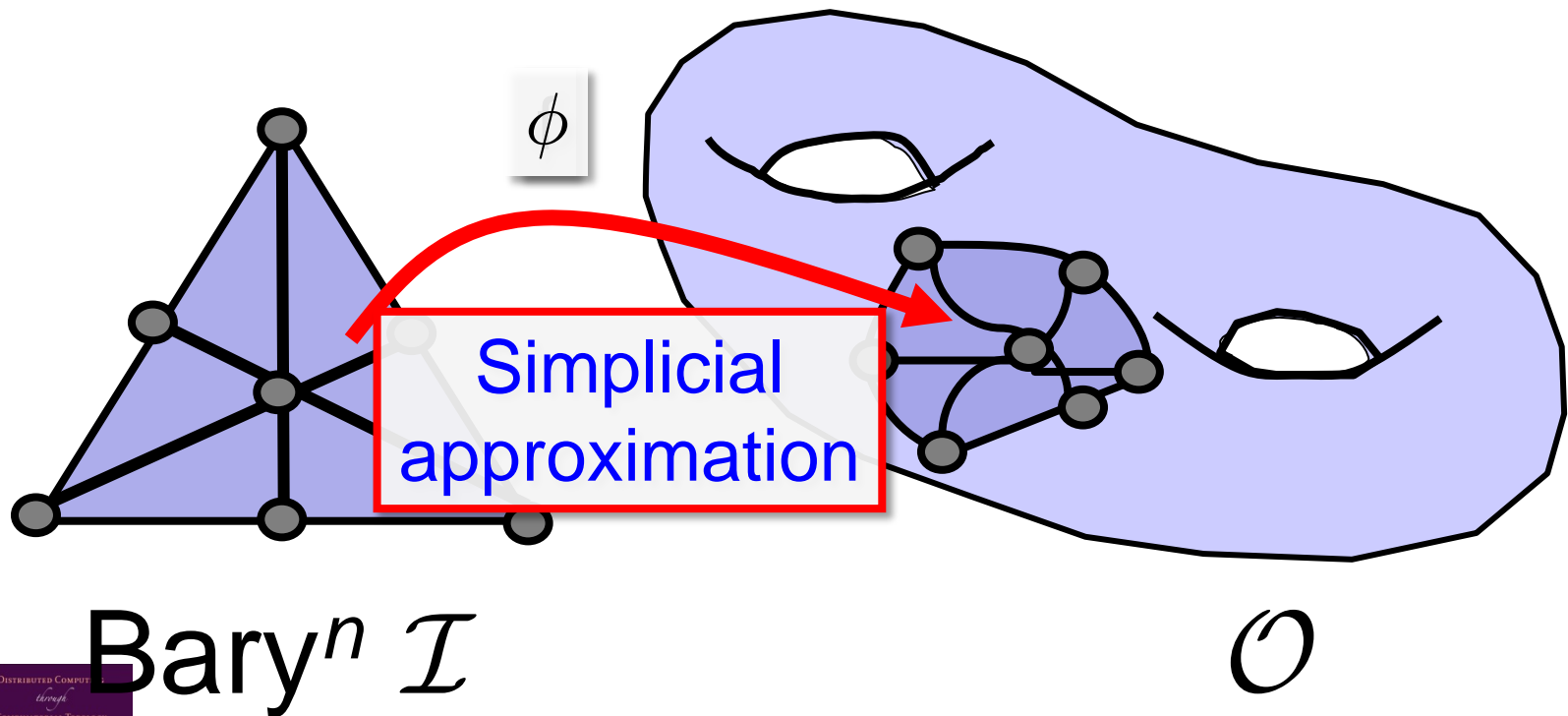


$Bary^n \mathcal{I}$

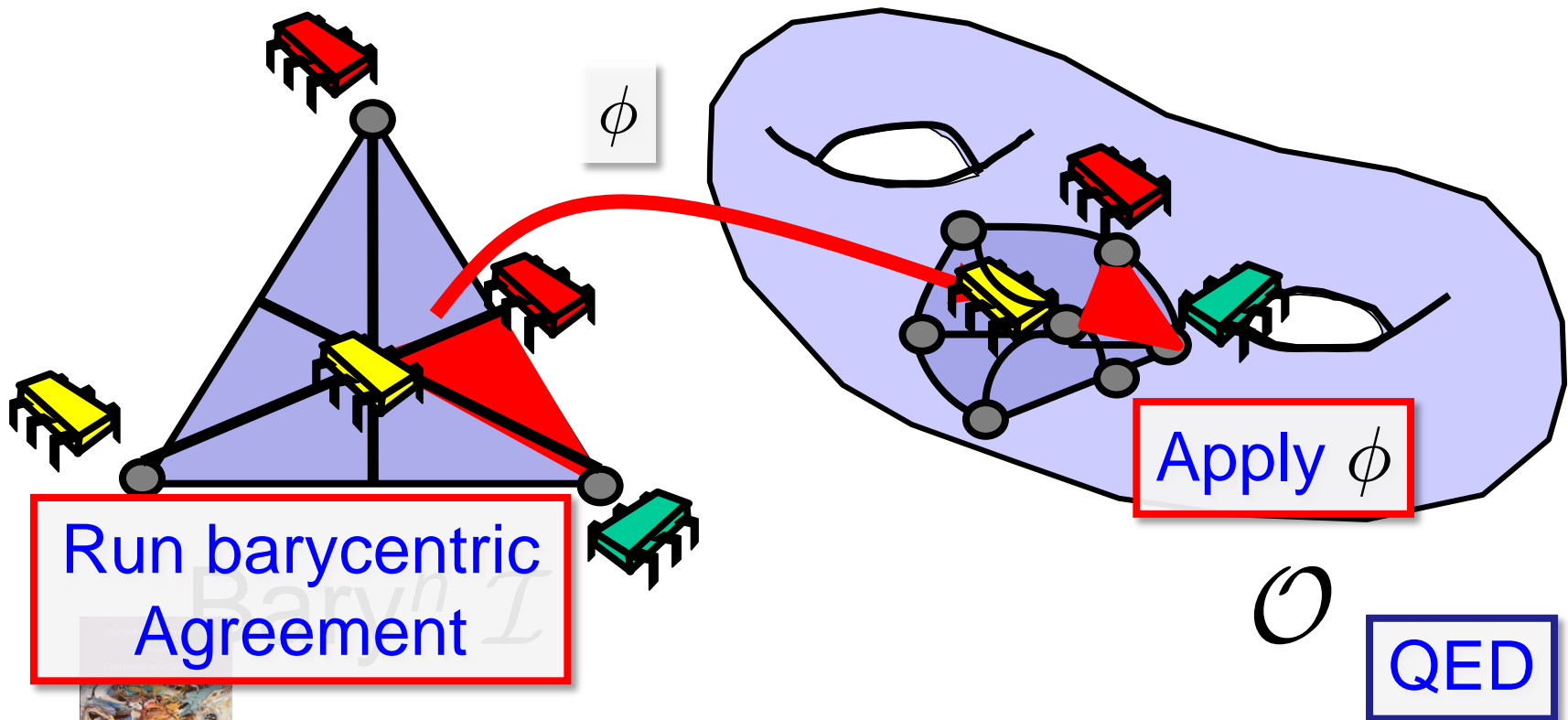
$\mathcal{O}$



# Map $\Rightarrow$ Protocol



# Map $\Rightarrow$ Protocol



# Protocol $\Rightarrow$ Map

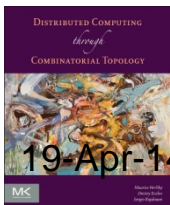
Lemma

If ...

there is a WF-RW protocol for  
 $(\mathcal{I}, \mathcal{O}, \Delta)$  ...

then ... only if ...

there is a continuous  
 $f: |\text{skel}^n \mathcal{I}| \rightarrow |\mathcal{O}|$  carried by  $\Delta$ .



19-Apr-14

# Protocol $\Rightarrow$ Map

Proof strategy

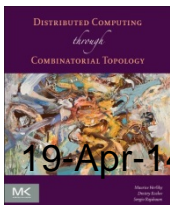
Inductive construction

$$g_d: |\text{skel}^d \mathcal{I}| \rightarrow |\Xi(\mathcal{I})|.$$

Base  $d = 0$

Define  $g_0: |\text{skel}^0 \mathcal{I}| \rightarrow |\Xi(\mathcal{I})| \dots$

Let  $g_0(v)$  be any vertex in  $\Xi(\{v\})$



19-Apr-14

# Protocol $\Rightarrow$ Map

Induction Hypothesis

$$g_{d-1}: |\text{skel}^{d-1} \mathcal{I}| \rightarrow |\Xi(\mathcal{I})|$$

For all  $\sigma$  in  $\text{skel}^{d-1} \mathcal{I}$

$$g_{d-1} \text{ sends } |\text{skel}^{d-1} \sigma| \mapsto |\Xi(\sigma)|$$

But  $\Xi(\sigma)$  is  $(d-1)$ -connected (earlier theorem)

Can extend to  $g_d: |\sigma| \mapsto |\Xi(\sigma)|$

Yielding  $g_d: |\text{skel}^d \mathcal{I}| \rightarrow |\Xi(\mathcal{I})|$

# Protocol $\Rightarrow$ Map

Constructed

$$g: |\text{skel}^n \mathcal{I}| \rightarrow |\Xi(\mathcal{I})|$$

Simplicial decision map

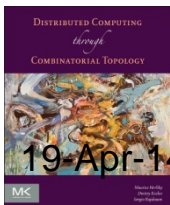
$$\delta: \Xi(\text{skel}^n \mathcal{I}) \rightarrow \mathcal{O}$$

$$|\delta|: |\Xi(\text{skel}^n \mathcal{I})| \rightarrow |\mathcal{O}|$$

Composition  $f = |\delta| \cdot g$  yields

$$f: |\text{skel}^n \mathcal{I}| \rightarrow |\mathcal{O}| \text{ carried by } \Delta.$$

QED



This work is licensed under a [Creative Commons Attribution-ShareAlike 2.5 License](https://creativecommons.org/licenses/by-sa/3.0/).

- **You are free:**
  - **to Share** — to copy, distribute and transmit the work
  - **to Remix** — to adapt the work
- **Under the following conditions:**
  - **Attribution.** You must attribute the work to “Distributed Computing through Combinatorial Topology” (but not in any way that suggests that the authors endorse you or your use of the work).
  - **Share Alike.** If you alter, transform, or build upon this work, you may distribute the resulting work only under the same, similar or a compatible license.
- For any reuse or distribution, you must make clear to others the license terms of this work. The best way to do this is with a link to
  - <http://creativecommons.org/licenses/by-sa/3.0/>.
- Any of the above conditions can be waived if you get permission from the copyright holder.
- Nothing in this license impairs or restricts the author's moral rights.

