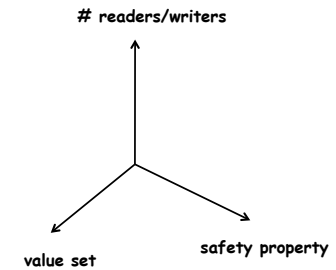


## FDS 2012: Atomic and immediate snapshots

© 2012 P. Kuznetsov

## The space of registers

- Nb of writers and readers: from 1W1R to NWNR
- Size of the value set: from binary to multi-valued
- Safety properties: safe, regular, atomic



All registers are (computationally) equivalent!



© 2012 P. Kuznetsov

2

## Transformations

From 1W1R binary safe to 1WNR multi-valued atomic

- I. From safe to regular (1W1R)
- II. From one-reader to multiple-reader (regular binary or multi-valued)
- III. From binary to multi-valued (1WNR regular)
- IV. From regular to atomic (1W1R)
- V. From 1W1R to 1WNR (multi-valued atomic)
- VI. From 1WNR to NWNR (multi-valued atomic)
- VII. From safe bit to atomic bit (optimal)



© 2012 P. Kuznetsov

3

## This class

- Atomic snapshot: reading multiple locations atomically
  - ✓ Write to one, read all
- Immediate snapshot
  - ✓ Write-read in groups
- Topological representation of shared-memory computing



4

## Atomic snapshot: sequential specification

- Each process  $p_i$  is provided with operations:
  - ✓  $\text{update}_i(v)$ , returns ok
  - ✓  $\text{snapshot}_i()$ , returns  $[v_1, \dots, v_N]$
- In a **sequential** execution:
  - For each  $[v_1, \dots, v_N]$  returned by  $\text{snapshot}_i()$ ,  $v_j$  ( $j=1, \dots, N$ ) is the argument of the last  $\text{update}_j(\cdot)$  (or the initial value if no such update)



© 2012 P. Kuznetsov

5

## Snapshot for free?

Code for process  $p_i$ :

**initially:**

shared 1W1R *atomic* register  $R_i := 0$

**upon  $\text{snapshot}()$**

```
 $[x_1, \dots, x_N] := \text{scan}(R_1, \dots, R_N)$  /*read  $R_1, \dots, R_N$ */  
return  $[x_1, \dots, x_N]$ 
```

**upon  $\text{update}_i(v)$**

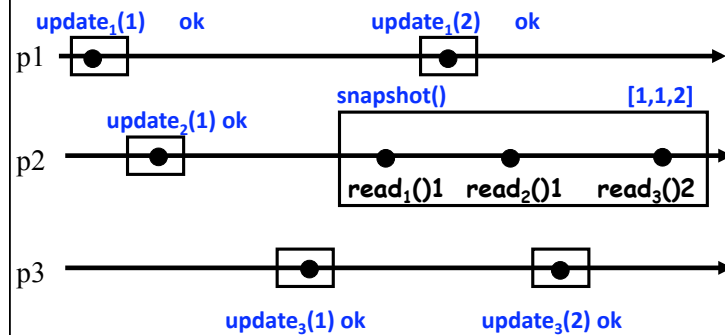
$R_i.\text{write}(v)$



© 2012 P. Kuznetsov

6

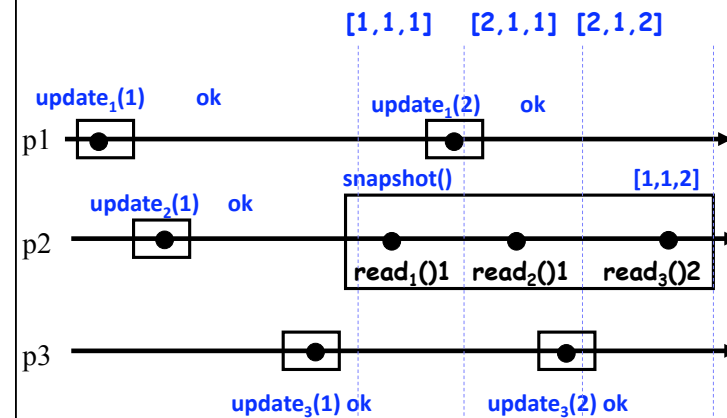
## Snapshot for free?



© 2012 P. Kuznetsov

7

## Snapshot for free?



© 2012 P. Kuznetsov

8

- What about 2 processes?
- What about **non-blocking** snapshot?
  - ✓ At least one correct process **makes progress** (completes infinitely many)



## Non-blocking snapshot

Code for process  $p_i$  (all written value are **unique**, e.g., equipped with a sequence number)

**Initially:**

shared 1W1R atomic register  $R_i := 0$

**upon snapshot()**

$[x_1, \dots, x_N] := \text{scan}(R_1, \dots, R_N)$

repeat

$[y_1, \dots, y_N] := [x_1, \dots, x_N]$

$[x_1, \dots, x_N] := \text{scan}(R_1, \dots, R_N)$

until  $[y_1, \dots, y_N] = [x_1, \dots, x_N]$

return  $[x_1, \dots, x_N]$

**upon update<sub>i</sub>(v)**

$R_i.\text{write}(v)$

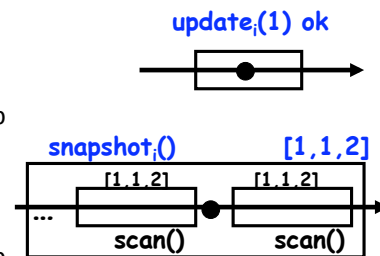


## Linearization

Assign a **linearization point** to each operation

- update<sub>i</sub>(v)
  - ✓  $R_i.\text{write}(v)$  if present
  - ✓ Otherwise remove the op
- snapshot<sub>i</sub>()
  - ✓ if complete – any point between identical scans
  - ✓ Otherwise remove the op

Build a **sequential history S** in the order of linearization points

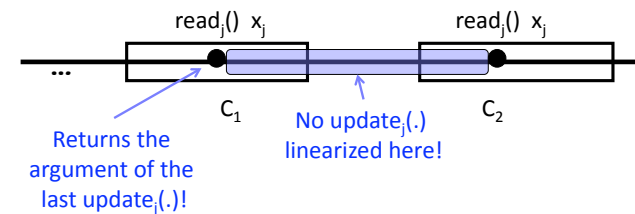


## Correctness: linearizability

S is legal: every snapshot<sub>i</sub>() returns the last written value for every  $p_j$

Suppose not: snapshot<sub>i</sub>() returns  $[x_1, \dots, x_N]$  and some  $x_j$  is **not** the argument of the last update<sub>j</sub>(v) in S preceding snapshot<sub>i</sub>()

Let  $C_1$  and  $C_2$  be two scans that returned  $[x_1, \dots, x_N]$



## Correctness: non-blockingness

An  $\text{update}_i()$  operation is wait-free (returns in a finite number of steps)

Suppose process  $p_i$  executing  $\text{snapshot}_i()$  eventually runs in isolation (no process takes steps concurrently)

- All scans received by  $p_i$  are distinct
- At least one process performs an update between
- There are only finitely many processes  $\Rightarrow$  at least one process executes infinitely many updates



© 2012 P. Kuznetsov

13

## General case: helping?

What if an update interferes with a snapshot?

- Make the update do the work!

**upon snapshot()**

```
[x1, ..., xN] := scan(R1, ..., RN)
[y1, ..., yN] := scan(R1, ..., RN)
if [y1, ..., yN] = [x1, ..., xN] then
    return [x1, ..., xN]
else
    let j be such that
        xj ≠ yj and xj = (u, U)
    return U
```

**upon update<sub>i</sub>(v)**

```
S := snapshot()
Ri.write(v, S)
```

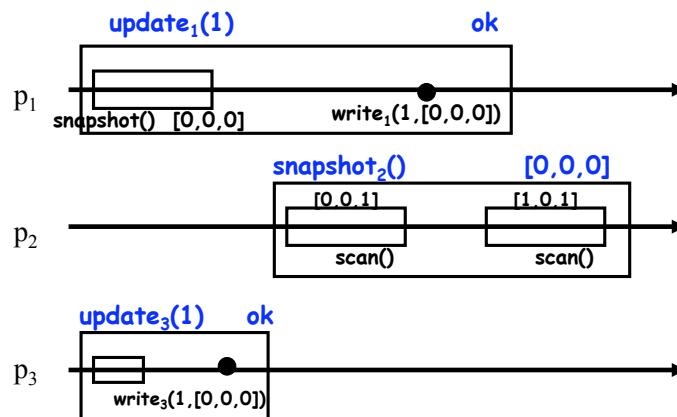
If two scans differ - some update succeeded to snapshot!  
Would this work?



© 2012 P. Kuznetsov

14

## Not that easy!



© 2012 P. Kuznetsov

15

## General case: wait-free atomic snapshot

**upon snapshot()**

```
[x1, ..., xN] := scan(R1, ..., RN)
while true do
    [y1, ..., yN] := [x1, ..., xN]
    [x1, ..., xN] := scan(R1, ..., RN)
    if [y1, ..., yN] = [x1, ..., xN] then
        return [x1, ..., xN]
    else if movedj and xj ≠ yj then
        let xj = (u, U)
        return U
    for each j: movedj := movedj ∨ xj ≠ yj
```

**upon update<sub>i</sub>(v)**

```
S := snapshot()
Ri.write(v, S)
```

If a process moved twice: its last snapshot is valid!



© 2012 P. Kuznetsov

16

## Correctness: wait-freedom

**Claim 1** Every operation (update or snapshot) returns in  $O(N^2)$  steps (bounded wait-freedom)

**snapshot:** does not return after a scan if a concurrent process moved and no process moved twice

- At most  $N-1$  concurrent processes, thus (pigeonhole), after  $N$  scans:
  - ✓ Either at least two consecutive identical scans
  - ✓ Or some process moved twice!

**update:** snapshot() + one more step



© 2012 P. Kuznetsov

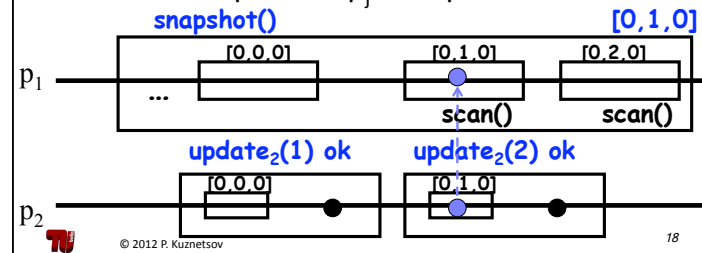
17

## Correctness: linearization points

**update<sub>i</sub>(v):** linearize at the  $R_i.write(v,S)$

complete **snapshot()**

- If two identical scans: between the scans
- Otherwise, if returned  $U$  of  $p_j$ : at the linearization point of  $p_j$ 's snapshot



© 2012 P. Kuznetsov

18

## The linearization is:

- Legal: every snapshot operation return the most recent value for each process
- Consistent with the real-time order: each linearization point is within the operation's interval
- Equivalent to the run (locally indistinguishable)

(Full proof in the lecture notes, Chapter 6)



© 2012 P. Kuznetsov

19

## One-shot atomic snapshot (AS)

Each process  $p_i$ :  
 $update_i(v_i)$   
 $S_i := snapshot()$

$S_i = S_i[1], \dots, S_i[N]$   
 (one position per process)

Vectors  $S_i$  satisfy:

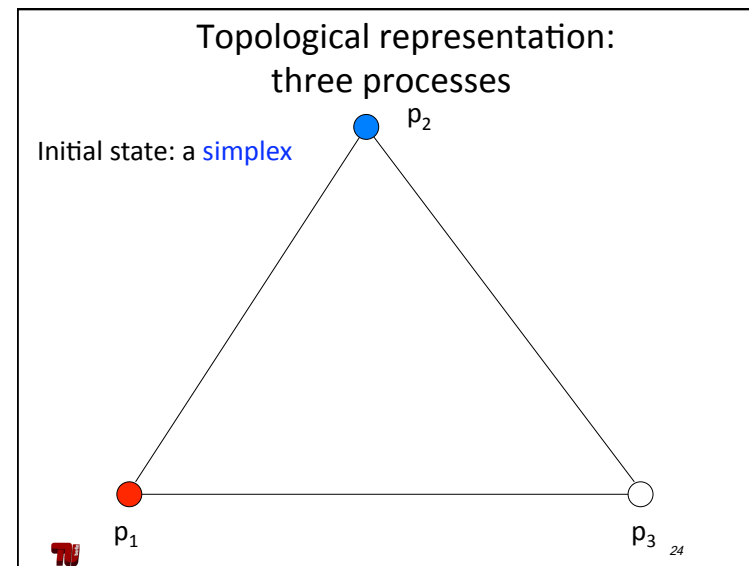
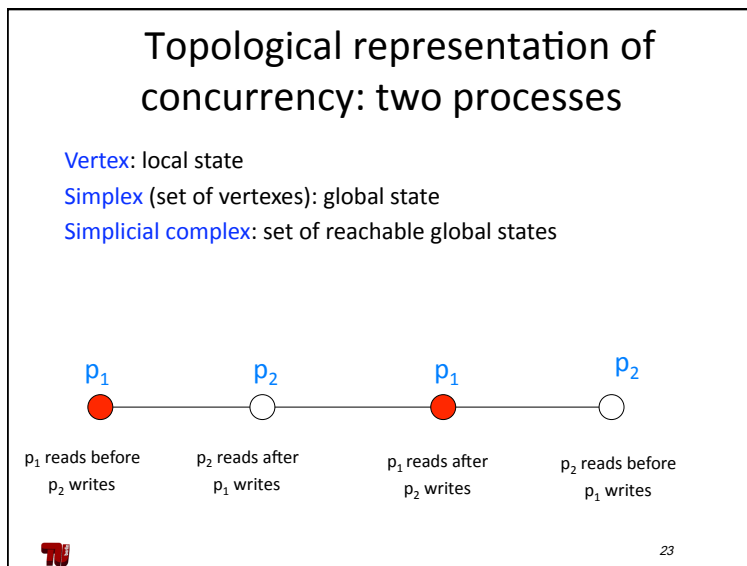
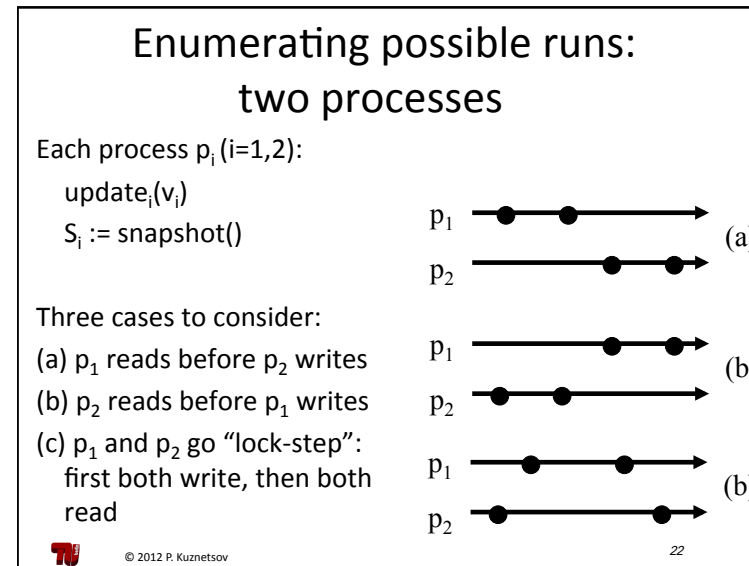
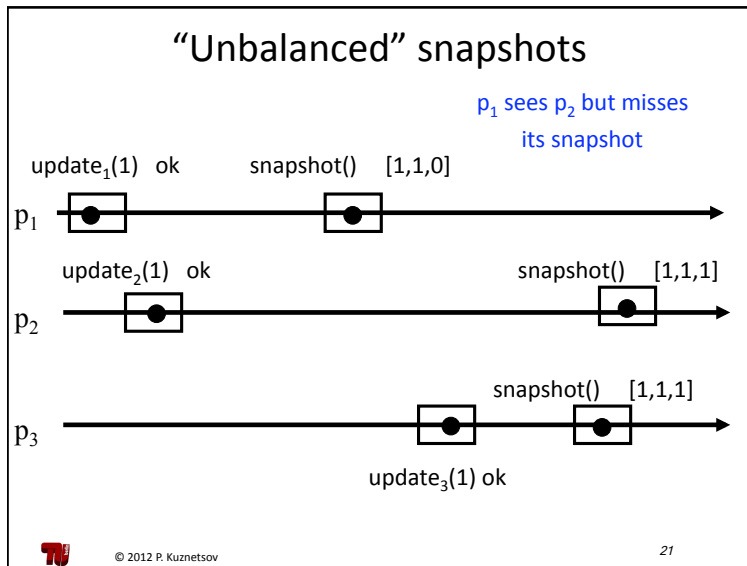
- Self-inclusion:** for all  $i$ :  $v_i$  is in  $S_i$
- Containment:** for all  $i$  and  $j$ :  $S_i$  is subset of  $S_j$  or  $S_j$  is subset of  $S_i$

Homework: prove the two properties are met by AS

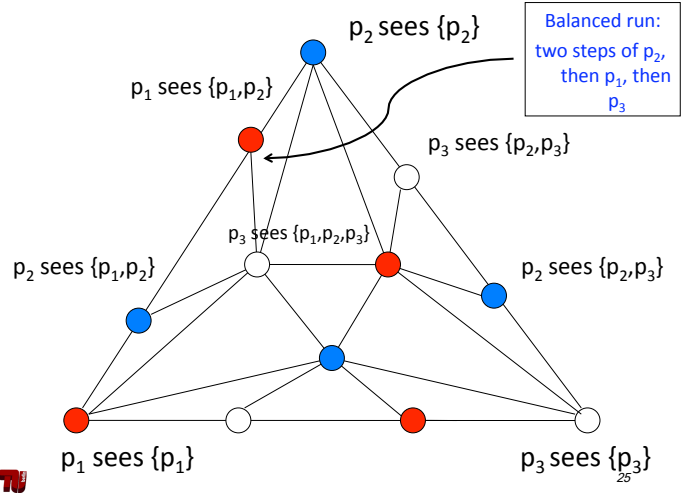


© 2012 P. Kuznetsov

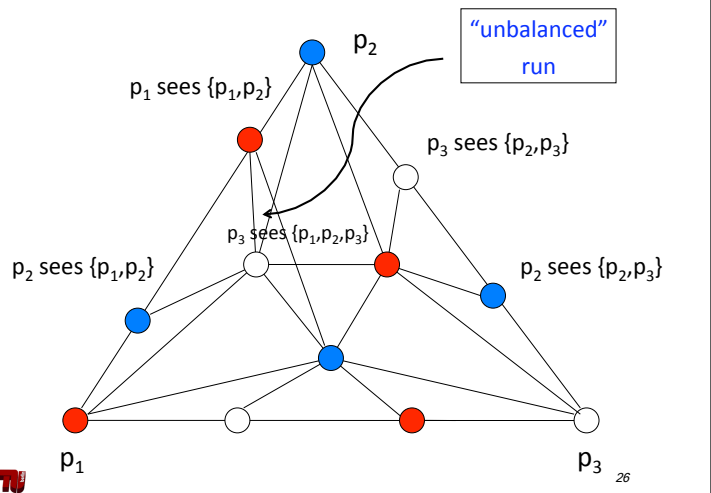
20



### Topological representation: one-shot AS



### Topological representation: one-shot AS



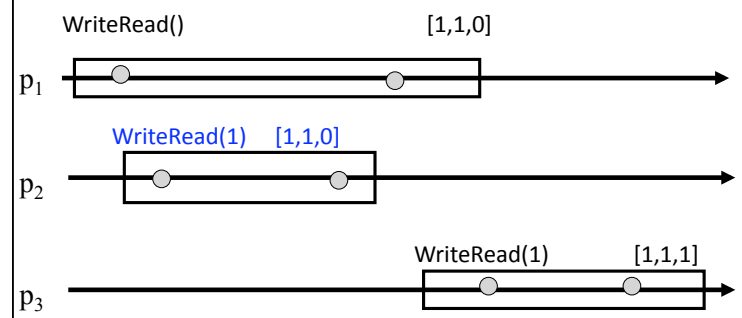
### One-shot *immediate* snapshot (IS)

One operation:  
WriteRead(v)

Each process  $p_i$ :  
 $S_i := \text{WriteRead}_i(v_i)$

- Vectors  $S_1, \dots, S_N$  satisfy:
- **Self-inclusion:** for all  $i$ :  $v_i$  is in  $S_i$
  - **Containment:** for all  $i$  and  $j$ :  $S_i$  is subset of  $S_j$  or  $S_j$  is subset of  $S_i$
  - **Immediacy:** for all  $i$  and  $j$ : if  $v_i$  is in  $S_j$ , then  $S_i$  is a subset of  $S_j$

### "Unbalanced" run eliminated



## Every IS run is “balanced”

The runs of IS match the **block** runs of one-shot AS that can be written as:

$$B_1 B_2 \dots B_k$$

where

- each  $B_i = \text{update}_{\Pi_i}() \text{ snapshot}_{\Pi_i}()$ 
  - ✓ All processes in  $\Pi_i$  update and then all processes in  $\Pi_i$  take snapshots
- $\Pi_1, \dots, \Pi_k$  is a partition of  $\{p_1, \dots, p_N\}$ 
  - ✓ E.g.,  $\Pi_1 = \{p_2\}$ ;  $\Pi_2 = \{p_1\}$ ;  $\Pi_3 = \{p_3\}$  in the balanced run above

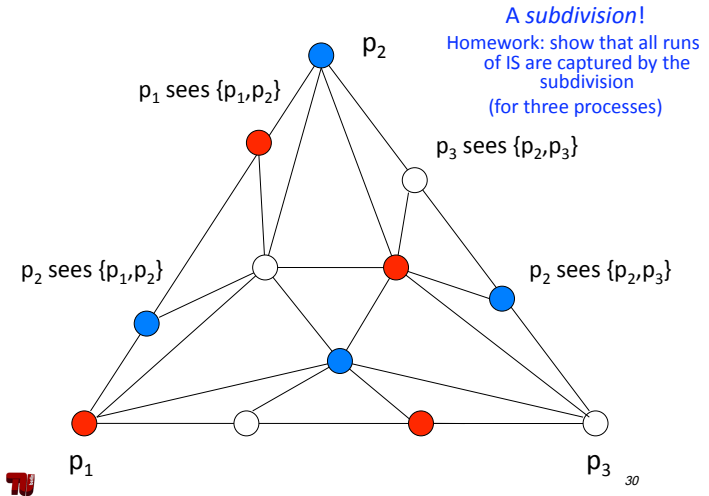
Homework: prove the equivalence of IS and block runs of AS



© 2012 P. Kuznetsov

29

## Topological representation: one-shot IS



## IS is equivalent to AS (one-shot)

- IS is a **restriction** of one-shot AS  $\Rightarrow$  IS is **stronger** than one-shot AS
  - ✓ Every run of IS is a run of one-shot AS
- Show that a few (one-shot) AS objects can be used to implement IS
  - ✓ One-shot ReadWrite() can be implemented using a series of update and snapshot operations



© 2012 P. Kuznetsov

31

## IS from AS

**shared variables:**

$A_1, \dots, A_N$  – atomic snapshot objects, initially  $[T, \dots, T]$

**upon ReadWrite<sub>i</sub>(v<sub>i</sub>)**

$r := N+1$

while true do

$r := r-1$  // drop to the lower level

$A_r.\text{update}_i(v_i)$

$S := A_r.\text{snapshot}()$

if  $|S|=r$  then //  $|S|$  is the number of non-T values in S

return S

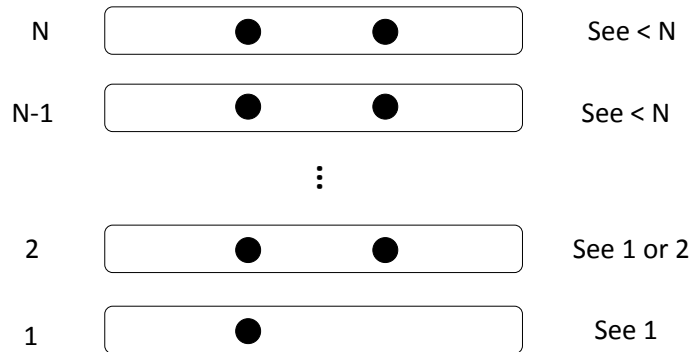


© 2012 P. Kuznetsov

32



## Drop levels



© 2012 P. Kuznetsov

33

## Correctness

The outcome of the algorithm satisfies Self-Inclusion, Snapshot, and Immediacy

- By induction on  $N$ : for all  $N > 1$ , if the algorithm is correct for  $N-1$ , then it is correct for  $N$
- Base case  $N=1$ : trivial



© 2012 P. Kuznetsov

34

## Correctness, contd.

- Suppose the algorithm is correct for  $N-1$  processes
- $N$  processes come to level  $N$ 
  - ✓ At most  $N-1$  go to level  $N-1$  and lower
  - ✓ (At least one process returns in level  $N$ )
  - ✓ Why?
- Self-inclusion, Containment and Immediacy hold for all processes that return in levels  $N-1$  or lower
- The processes returning at level  $N$  return **all  $N$  values**
  - ✓ The properties hold for all  $N$  processes! Why?



© 2012 P. Kuznetsov

35

## Iterated Immediate Snapshot (IIS)

Shared variables:

$IS_1, IS_2, IS_3, \dots$  // a series of one-shot IS

Each process  $p_i$  with input  $v_i$ :

$r := 0$

while true do

$r := r+1$

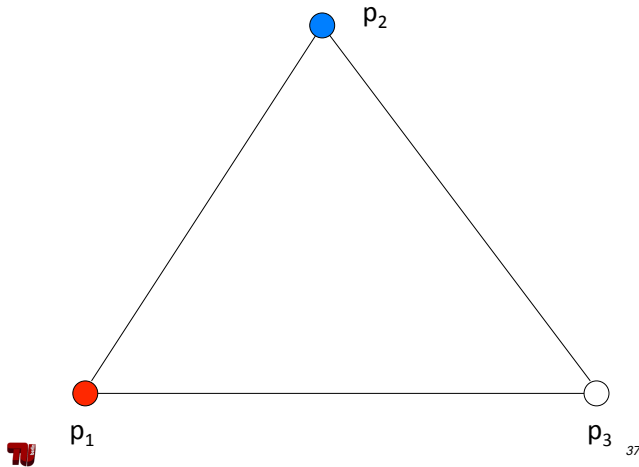
$v_i := IS_r.\text{WriteRead}_i(v_i)$



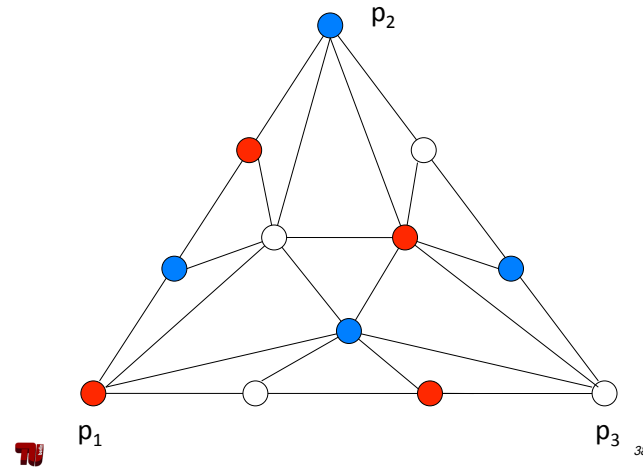
© 2012 P. Kuznetsov

36

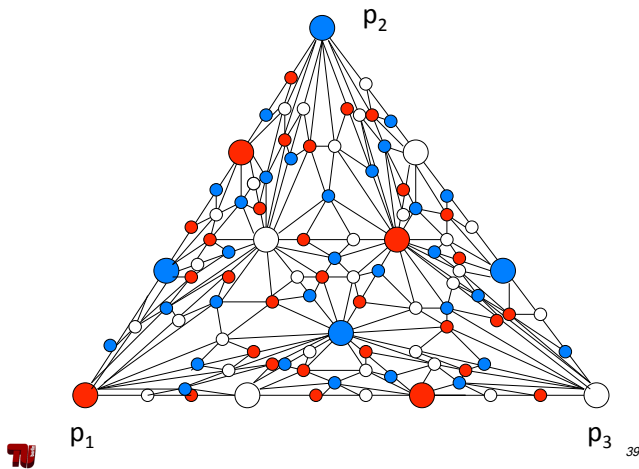
### Iterated standard chromatic subdivision



### ISDS: one round



### ISDS: two rounds



### IIS is equivalent to (multi-shot) AS

- AS can be used to implement IIS (wait-free)
  - ✓ Multiple instances of the construction above (one per iteration)
- IIS can be used to implement AS in the non-blocking manner:
  - ✓ At least one correct process performs infinitely many read or write operations
  - ✓ Good enough for protocols solving distributed tasks!