

CSC241S Advanced Distributed Computing

- Deterministic abortable (DA) obj.
 - wait free
 - linearizable

Like ordinary obj except ops that experience contention (i.e. are concurrent with other operations) may return "abort" (L) without having any effect on the object.

→ regular obj cannot be impl. with registers - need better sync primitives, so weaken.

→ background

- Obstruction free objects

Weakening of wait-free obj.

def obstruction-freedom = op that eventually executes solo (no contention) terminates and returns normal response -

fact $\forall n, \forall$ linearizable obs-free obj shared by n procs can be impl. using only registers -

exercise give an obs-free impl. of consensus for 2 procs using only registers -

• "Pausable" objects

- Every op invoked by correct process eventually returns ctrl to caller
- Op. that encounters no contention returns "normal" response -
- An operation that encounters contention can:
 - (a) return normal response
 - (b) ——— \perp and have no effect -
 - (c) ——— "pause" in which case it may or may not have taken effect - ~~At~~
 - at this point, caller must resume until (a) or (b).

fact $\forall n, \forall$ linearizable pausable obj shared by n procs can be impl. using only registers -

fact if pause is not an option, then \exists obj. that can't be impl. using only registers -

→ specifically abortable consensus (with ^{registers} pause) for 2 procs cannot be impl. using only

→ Wait-free consensus for 2 procs using a "DAC" obj. for 2 procs.

$D = 2$ -DAC objects

$X =$ registers

proof

$P_0: \text{propose}(v)$ $d_0 := D.\text{prop}(v)$ if $d_0 = \perp$ then return d_0 else $\underline{d_0}$	$P_1: \text{propose}(v)$ $rc := v$ repeat $d_i := D.\text{prop}$ until $d_i \neq \perp$ return d_i
---	--

can be impl using registers

→ Now def abortable objects
 Same as passable but (a) return normal resp.
 (b) return \perp in which case it may or may not have taken effect.

→ Query abortable objects
 Try to find out fate of aborted op -

• Back to DA objects -
 \perp : means did not take effect -

Every Def every obj has a type:

⚠ not allow 2 procs to access same point ⚠

$T = (OP, RES, Q, n, \delta)$

ops response states procs trans. func.

$\delta: Q \times OP \times \{1..n\} \rightarrow Q \times RES$
 $(q, op, i) \mapsto (q', res)$

DA counterpart of T:

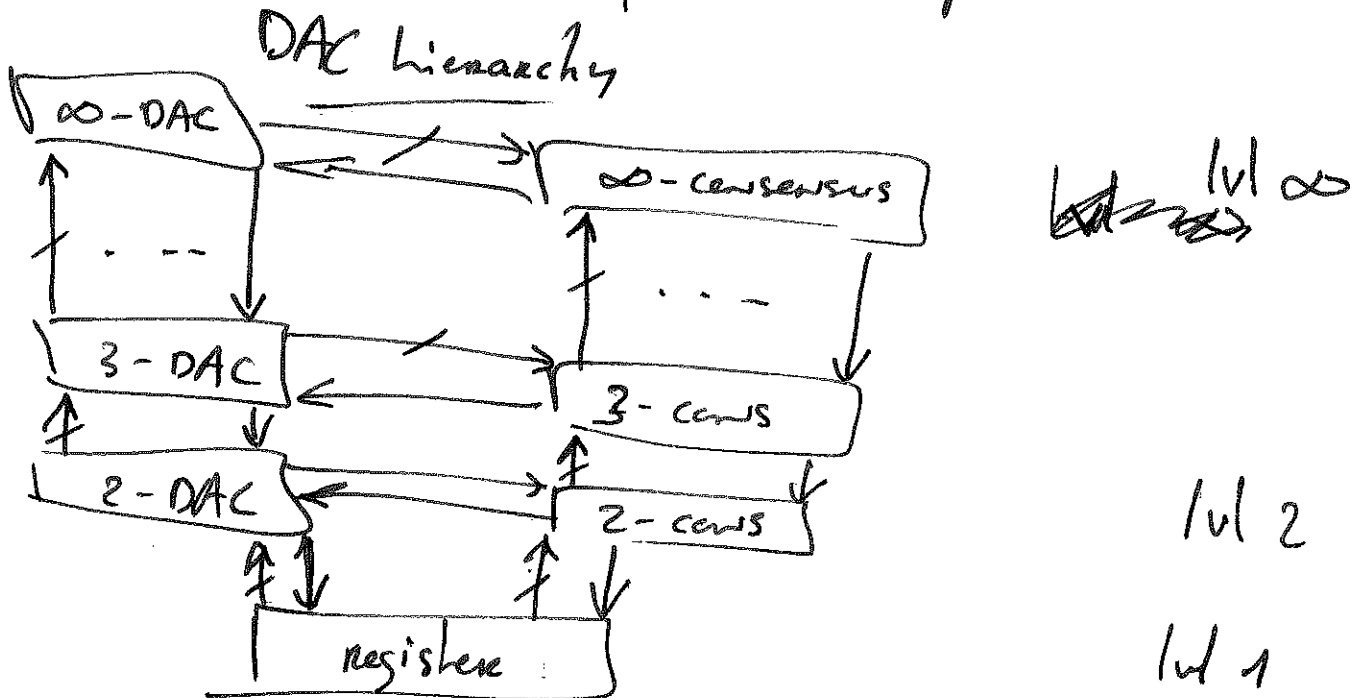
$$T^{da} = (OP, RES \cup \{\perp\}, Q, n, \int^{da})$$



def n -DAC : n -ported DA counterpart of n -consensus

Non-triviality (1) An op. can return \perp only if it is [safety] concurrent w/ another op.

(2) An op. that is interrupted by the crash [liveness] of the invoker can cause only a finite number of concurrent ops to abort.



Level^m of a T^{da} in DAC hierarchy = max n for which T^{da} & registers can impl. n -DAC

fact ~~DA~~ T & S^{da} ~~is~~ \rightarrow 1 registers \rightarrow 2-DAC
 $\cdot T$ & S^{da} ~~is~~ \rightarrow 1 registers \rightarrow 3-DAC

proof Sup. T & $S^{da} \rightarrow$ 3-DAC

then T & $S \rightarrow$ 3-DAC

\Rightarrow 2- T & $S \rightarrow$ 3-DAC

\Rightarrow 2-~~cons~~ \rightarrow 3-DAC

\Rightarrow 2-DAC \rightarrow 3-DAC

Conjecture If T is an ordinary type at level n of the Consensus hierarchy, then T^{da} is at level n of the DAC-hierarchy.

disproof ~~is~~

def sticky registers $\hat{=}$

initially Δ , read but write once -

idea DA-stick. reg.

FLP impossibility \rightarrow FD model

FD model step

- receive a msg
- query FD
- send a msg
- MR alg ($t < \frac{n}{2}$)
 - coordination at round r is $P_{R \bmod n}$ majority
 - decide iff all received msg from correct ones same.

Too general:

e.g. if $\exists p$ has initial val. 0 then ret \perp
else ret \emptyset

Not general enough:

e.g. Heartbeat FD

Failure Pattern:

Π : processes; \mathbb{N} = real time

$F: \mathbb{N} \rightarrow \mathcal{P}(\Pi)$

$F(t) = \{p \in \Pi \mid p \text{ crashed at time } t\}$

~~\mathbb{A}~~ Failure History (with range R) $R \subseteq \mathcal{P}(\Pi)$

$H: \mathbb{A} \times \mathbb{N} \rightarrow R$

$H(p, t) = \text{value ret by } p \text{'s FD module at } t.$

FD D : failure pattern \mapsto set of FD histories.

$D(F) =$ set of possible beh of D when Fail. Pt. is F .

Redefinitions:

$\diamond S$: $H \in \diamond S(F) \Leftrightarrow$ ^① $\{ \forall p \forall q \exists t, p \text{ connect}(F) \wedge q \notin \text{connect}(F) \wedge \forall \geq t \wedge q \in \mathcal{M}(p, t) \}$
 ② ~~then~~ p - then property.

Redef step:

formally: (p, m, d) ^{process} _{msg received} ^{FD value}

Redef schedule:

seq of steps.

Redef run of alg A using FD D .

$R_A = (F, \mathcal{M}, I, S, T)$ ^{seq of lines}
_{failure pat.} _{FD hist.} _{initial config} _{seq of steps in run (std)}
 $\rightarrow H \in D(F)$

props $\rightarrow \forall t \in T$, if p crashed at t' ~~then~~ $t < t'$

Environment def:

$E =$ set of failure pat.

Alg A solve P using FD D in env E

• Weakest FD to solve P in env \mathcal{E} .

• Def: $D \succeq_{\mathcal{E}} D'$ (\Leftrightarrow) D provides at least as much info as D' (in \mathcal{E})
 $\Leftrightarrow \exists$ alg that uses D to simulate D' (in \mathcal{E}).

e.g. $P \succeq_{\mathcal{E}} HB$ (change set of procs into #)

• Leader FD Ω

• Quorum PD Σ

• output = quorum: set of procs

• reqs: (1) any 2 quorums (at any times and any processes) must intersect.

(2) eventually quorum of correct process must contain only correct -

• Weakest FD to solve consensus in env. \mathcal{E} .

• FD $\chi_{\mathcal{E}}$ s.t. (1) \exists alg that uses $\chi_{\mathcal{E}}$ to solve consensus.

(2) \forall FD D that ~~uses $\chi_{\mathcal{E}}$ to solve~~ can be used to solve consensus,

$$D \succeq_{\mathcal{E}} \chi_{\mathcal{E}}$$

• Thm $\forall \mathcal{E}$, the weakest FD to solve consensus is $\left(\begin{array}{c} \text{anecdote} \\ \Omega, \Sigma \end{array} \right)$ ^{leader quorum}

MR $\Leftrightarrow (\Omega, \Sigma)$ intersect^o prop gives maj facon quorum.

Ben-OR also can be

For (1): modify MR

(2): $\forall \mathcal{E}, \forall D$ that can be used to solve consensus ^{in \mathcal{E}}

(a) $D \succeq_{\mathcal{E}} \Omega$ [CHT 92]

(b) $D \succeq_{\mathcal{E}} \Sigma$ [DF03]

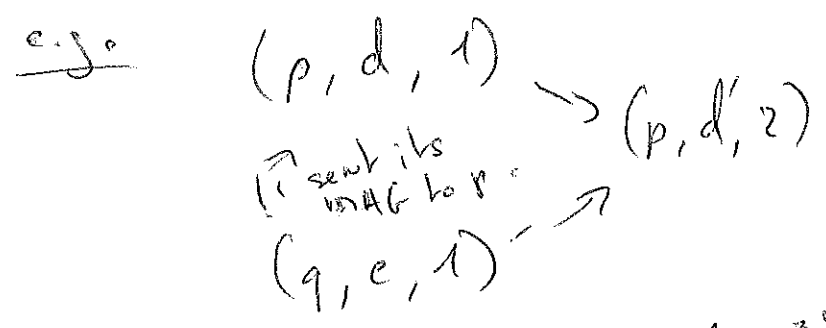
Exercise: implement Σ in a synch if $E \rightarrow \mathbb{R}^+ \tau < \frac{n}{2}$
 → wait for maj of

Proof of (b)
 Any FD solving consensus in E can be transformed to D^* in E wlog

Let alg be any alg solving consensus in E .

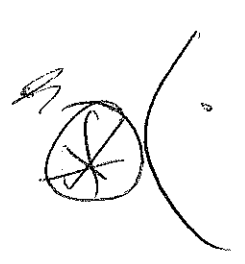
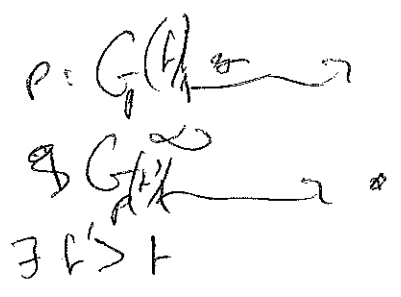
Two interlocking components: directed acyclic graph

- (1) DAG building
- (2) extraction component (of Σ)



fair if connected, has infinitely many samples.

- ① G_p monotonically increases $\Rightarrow G_p \xrightarrow{\infty}$ limit graph \downarrow of P
- ② If $(q, d, k) \rightarrow (q', d', k')$ in G_p then q took k 'th sample and said d before q' took its k' 'th sample and said d'
- ③ if p is connect $G_p \xrightarrow{\infty}$ contains a "fair" path (all connect \mathbb{R}^+ have only many samples on that path)



Simulated schedules of A .
 run $R = (P, A, -, -, -)$ of DAG-bldg \downarrow of P

e.g. Fix g , path in G_p .

$$g: (p_1, d_1, k_1) \rightarrow (p_2, d_2, k_2) \rightarrow (p_3, d_3, k_3)$$

note: no need to be consecutive k s,
even for same process ($\sigma \rightarrow \sigma$?)

- Fix initial config I of A
- Let S : set of schedules of A that are:
 - applicable to I
 - compatible with g



$$S = (p_1, m_1, d_1) \quad (p_2, m_2, d_2)$$

$\neq \text{null}$

Simultaneous schedules of A

consider $R = (F, H, \text{star})$

g : path in some DAG

I : init config of A

G_p | (S1) If S is comp w/ g and applicable to I then $\exists T$, s.t. (F, H, I, ST) is a run of A .

G_p^∞ | (S2) If g is fair then $\exists S^\infty$ that is compatible w/ g & appl. to I and $\exists T^\infty$, $(F, H, I, S^\infty, T^\infty)$ is an admissible run of A .

~~Assume~~ $Sch(G, I) =$ set of schedules compat w/ some path $\in G$ and appl. to I .

- p maintains variable $u_p =$ "recent" sample taken by p .
- p keeps building its DAC while

G_p restricted to nodes u_p and taken' $Sch(G_p | u_p, I_0)$ contains S_0 such that p decides in $S_0(I)$. all inputs = 0

and $Sch(G_p | u_p, \bar{I}_1)$ contains $S_1(I_1)$ all inputs = 1

When this happens, p computes new quorum

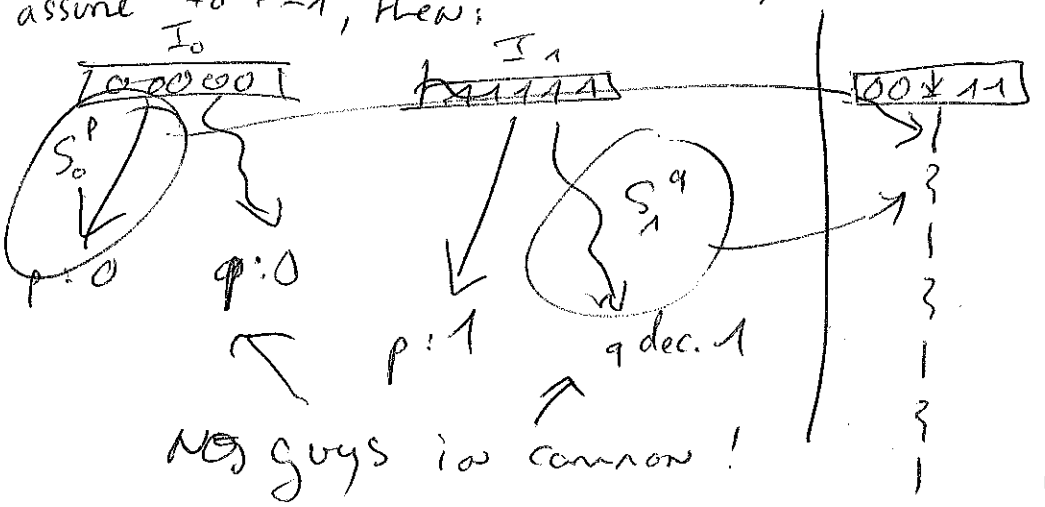
ok. $\Sigma\text{-out}_p = \text{participants}(S_0) \cup \text{part.}(S_1)$
 $u_p :=$ most recent sample of p .

careful about things happening at the same time -

~~$Sch(G, I)$~~
 proof of prop of quorum

~ actually continue pointing to things about?

assume $I_0 \neq I_1$, then:



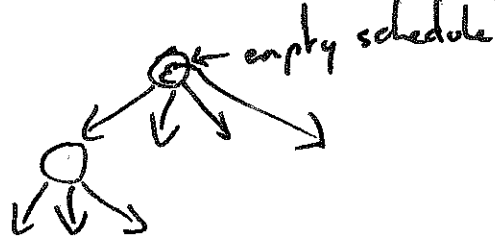
construct a 3rd run from 2 possible runs, then we don't have consensus!

A solves consensus using D in E

→ show $D \geq_E \Omega$

Organise initial configs $I_0 \dots I_m$ ($2^{\text{processes}}$) in tree

→ Simulation Tree T_G^I

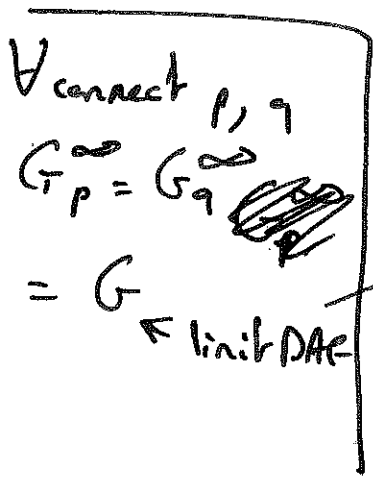


nodes: scheds in $Sch(G, I_i)$

edges: $S \rightarrow S'$ iff $S' = S \dot{\cup} \langle p, m, d \rangle$

→ Simulation Forest \tilde{F}_G

$$\tilde{F}_G = \{T_G^{I_0}, T_G^{I_1}, \dots, T_G^{I_m}\}$$

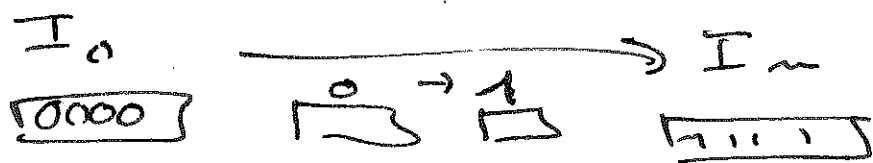


\tilde{F}_G has infinite length path (if sched.)
+ fan out (for FD infinite domain)

→ Tagging nodes in T_G^I

tag node v of T_G^I with $0, 1, \{0, 1\}$
unival.
bivalent

Fact 1: \tilde{F}_G contains a tree w/ bivalent root or $\exists i, 0 \leq i < m$ s.t. root of $T_G^{I_{i-1}}$ is 0-val & root of $T_G^{I_i}$ is 1-val.



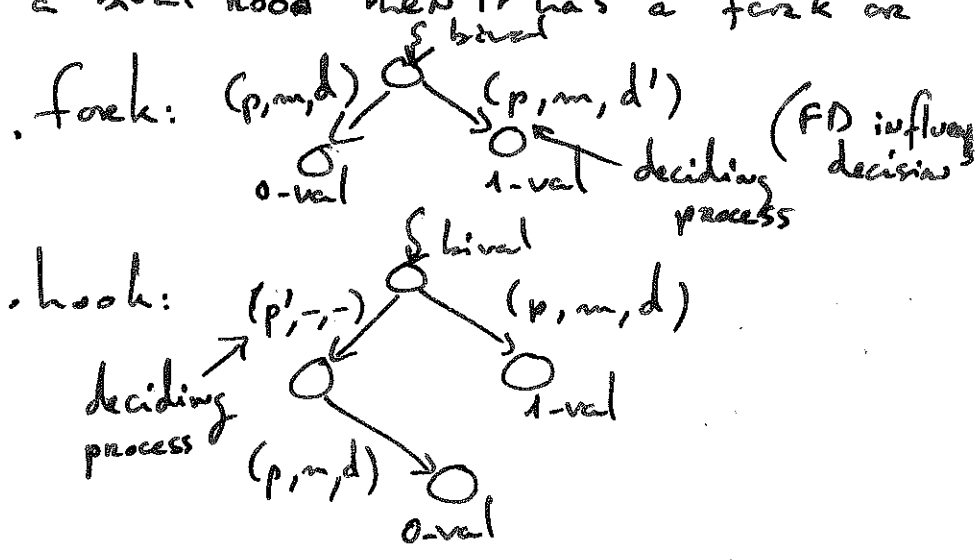
Fact 2: If root of $T_G^{I_{i-1}}$ is 0-val & root of T_G^I is 1-val then p_i is correct.

proof: by contradiction, p_i changes the decision.

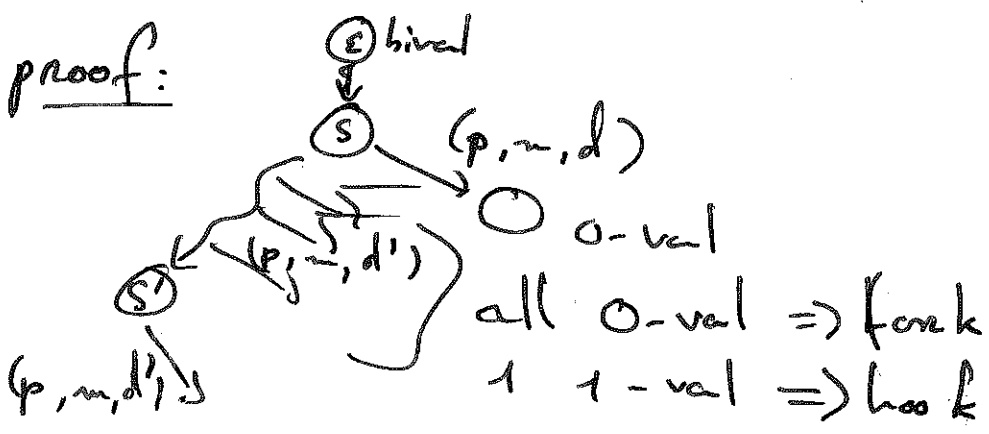
Fact 3: If a tree has a bival root then it has a bival node s.t. $\exists p, \exists m, S \cdot (p, m, d)$ is a node and

\forall node of the form $S \cdot E \cdot (p, m, -)$ is unival.
proof: by contradiction. show break consensus

Fact 4: If a tree has a bival root then it has a fork or a hook:

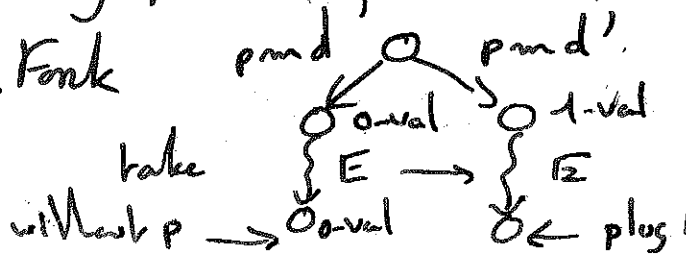


proof:



Fact 5: The deciding process of a hook / fork is correct.

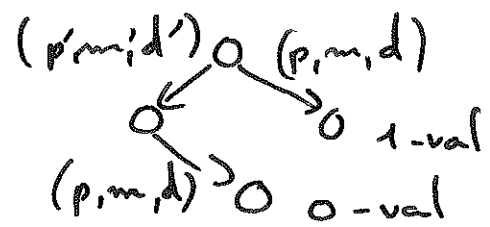
proof: Fork



Assume p is \rightarrow correct then $\exists z$ such that compatible w/ right side.

Hook:

case 1: $p = p'$
(same as fork)



Comment: FD is not nonotone!
 $(p, s, t_1) \rightarrow$
 $(p, t, t_2) \rightarrow$

case 2, $p \neq p'$

E w/out p' ~~and~~ taking steps after d' .

-> Rule for picking connect p from \mathcal{F}_G

Let i be min ~~of~~ $0 \leq i \leq n$, s.t. either (a) $T_G^{I_i}$ is bival.
 or (b) $T_G^{I_{i-1}}$ is 0-val and $T_G^{I_i}$ is 1-val. If (b), then
 pick p_i is connect (by fact 2). If (a), then choose smallest
 hook or fork (in some encoding of graphs) and pick the
 deciding process (by fact 5).

-> Applying rule to process:

Each p based on its DAG G_p constructs its own simulation forest.

Fact 6: For $\forall S \in \mathcal{F}_G, \exists t, \forall$ connect $p, \forall t' \geq t, \mathcal{F}_p(t')$
 contains S & connect tags for S .

Have to use fact 6, bc $\exists t, G_p(t) = G_q(t)$ ^{for} _{connect} _{p, q}
 is not true, since it is possible for 2 connect
 processes one constantly exchanging messages -

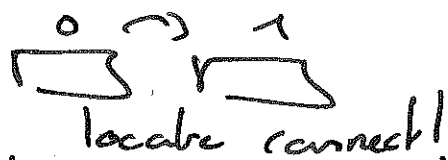
So now, just wait for fact 6, then that all tags have
 been applied (s.t. all connect have ^{same} min. fork/hook) -

FLP \rightarrow Lemma 1: 3 bivalency config -
 (origin of proof)

Lemma 2: 3 lock -

A string
 consensus w/ 0

\rightarrow if it violates lemma 1 then



2 then hook
 the guy who takes step (decides) is correct!

\rightarrow How to extract Ω when A solves non-unif consensus

\rightarrow how to tag? cannot know correct. \leftarrow

NU agreement. not 2 correct
 poss decide \neq

just tag iff in subtree has self deciding tag

history: $q \rightarrow p$ good hypo 1, 1

\rightarrow How to extract Σ

\rightarrow cannot use same argus before, but can extract Σ^0 (where interesec^o property holds for correct)

\rightarrow Weakest FD to solve NU consensus:

(Ω, Σ^0) .

— End of class —

1. Cover Welsh

2. Present set-binariness (k-set agreement)

3. Ω_j with ~~tags~~