

DISTRIBUTED SYSTEMS

LECTURE 2

WHERE WE ARE

- WHAT? ALGORITHMS/PROTOCOLS THAT RUN ON MULTIPLE NODES
 - ↳ CAN COMMUNICATE OVER A NETWORK
- FOCUS ON MESSAGE PASSING
 - ↳ NODES COMMUNICATE BY SENDING & RECEIVING MESSAGES
- ASSUME THE ASYNCHRONOUS MODEL
 - ↳ MESSAGES CAN BE DROPPED OR ARBITRARILY DELAYED
 - FAIRNESS
- WHY DO ANY OF THIS?
 - FAULT TOLERANCE

- SCALING

- CHALLENGE: DISTINGUISHING BETWEEN **FAILED** & **SLOW** NODES

↳ NODE MIGHT BE SLOW OR MESSAGES DELAYED

TODAY

- TOOLS FOR REASONING ABOUT DISTRIBUTED ALGORITHMS & PROTOCOLS

→ EVENTS & HOW TO ORDER THEM

→ TRACES

→ CORRECTNESS PROPERTIES: SAFETY & LIVENESS

- A 'SIMPLE' PROTOCOL FOR CAPTURING THE STATE OF A DISTRIBUTED SYSTEM



To check correctness

SPECIFYING ALGORITHMS

$x, y \in \mathbb{Z}$

SWAP (x, y) :

1 $x = x \wedge y$

2 $y = x \wedge y$

3 $x = x \wedge y$

$x = y_{in}; y = x_{in}$

1 } correct?
3 }
2 }

... thinking about algorithms as sequential

Generally used to thinking about algorithms = steps.

FOR LONG RUNNING SYSTEMS, ALSO COMMON TO THINK ABOUT CORRECTNESS IN TERMS OF ORDER IN WHICH THINGS OCCUR

EG. LINEARIZABILITY (2 WEEKS FROM NOW)

$WRITE(x, 1) \prec_{HB} READ(x)$

$\Rightarrow READ(x)$ RETURNS 1 OR NEWER WRITE

TODAY'S QUESTION:

IN A DISTRIBUTED SYSTEM
HOW DO WE DECIDE WHETHER

EVENT e_0 HAPPENS-BEFORE e_1 ?

EVENTS: THINGS A PROCESS DOES

FROM LAST CLASS

ON RECV(m) FROM X:
 COMPUTE + UPDATES

SEND(y, m₀)
SEND(z, m₁)
 ...

- RECEIVES
- SENDS
- WHY NOT COMPUTE?

ENVIRONMENTAL

- PROCESS FAILURES
 - PROCESS LAUNCHES
- MIGHT NOT BE UNDER THE PROGRAM'S CONTROL

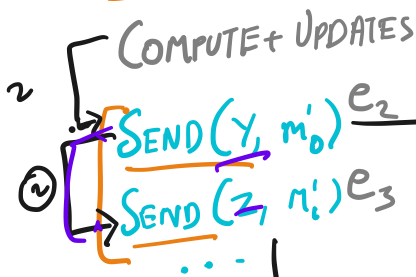
send(p, m₁)
send(q, m₂)

send(p, m₁)
send(q, m₂)

CAUSALITY

SOME EVENTS HAVE A NATURAL NOTION OF ORDERING

ON RECV(m) FROM X: e₁



① MESSAGE MUST BE SENT BEFORE IT IS RECEIVED

ON RECV(m₀) FROM... : e₄

ON RECV(m₁) FROM... : e₅

② ASSUME THAT EACH PROCESS EXECUTES PROGRAM SEQUENTIALLY

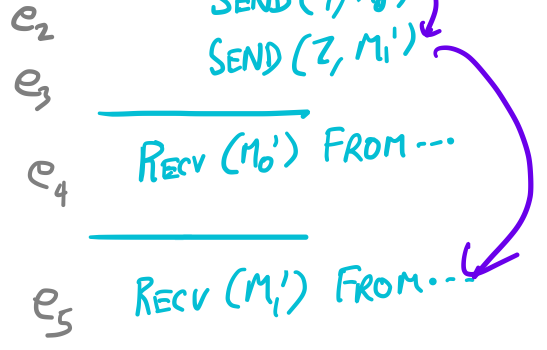
③ TRANSITIVE CLOSURE OF RELATIONS FROM 1 & 2.

PARTIAL ORDER

e₁ RECV(m) FROM X
 SEND(y, m₀)

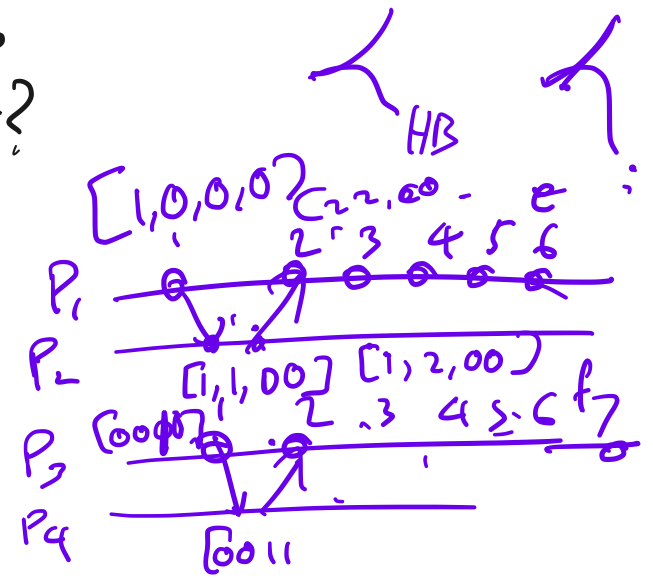
e₁ e₄?

|||

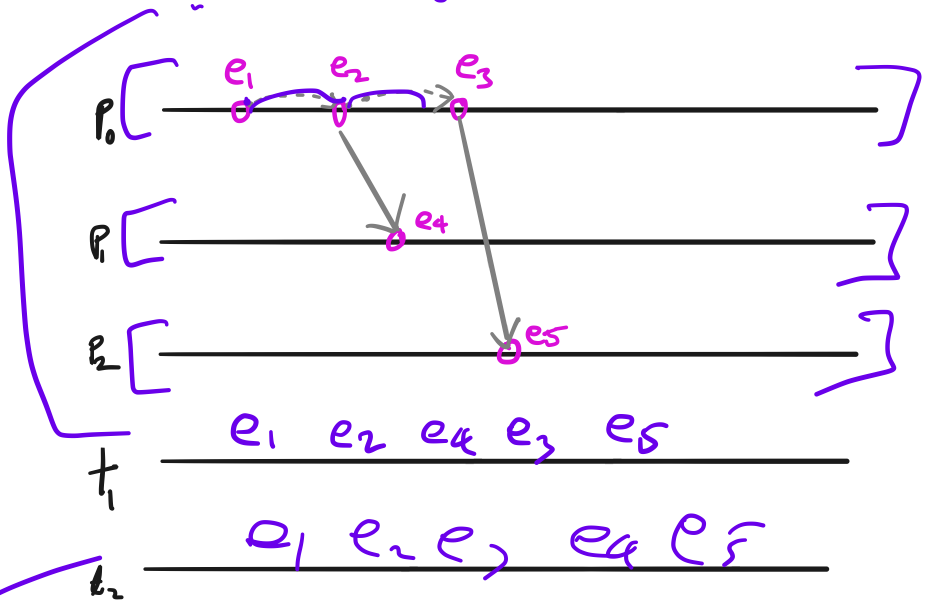
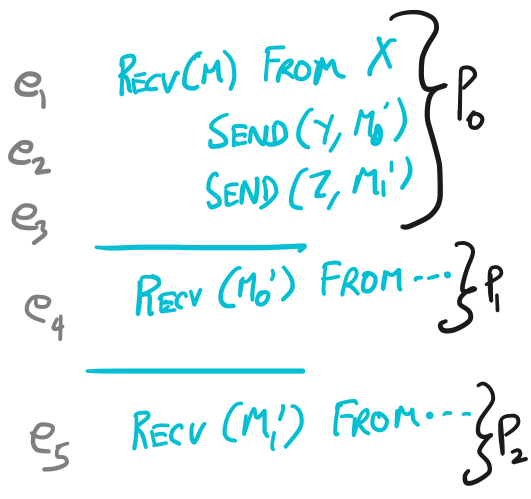


e_4 e_5 ?
 e_3 e_4 ?
 e_2 e_5 ?

A.S.



PRODUCING A TRACE
 ↳ Total order

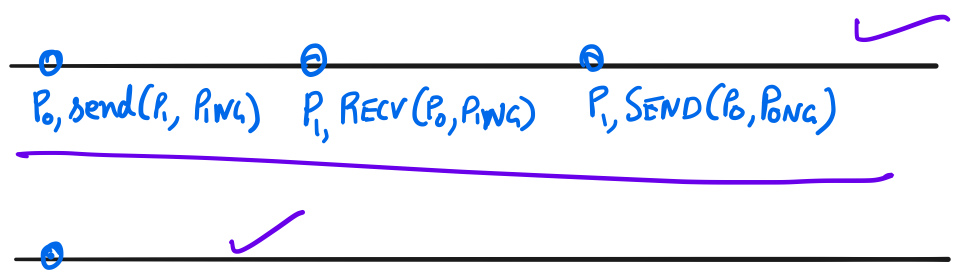


TRACE

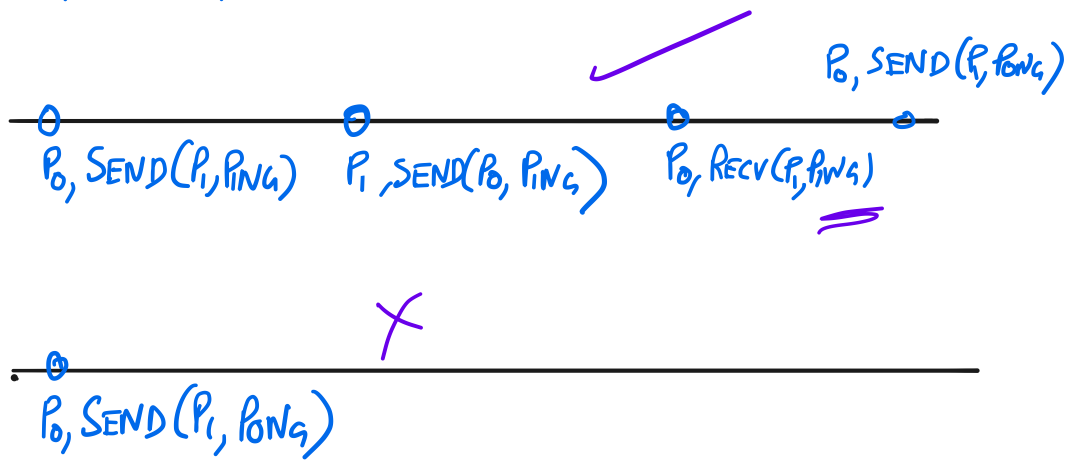
CORRECTNESS CONDITIONS FOR ADISTRIBUTED PROTOCOL (OR ALGORITHM) LIMIT WHAT

TRACES ARE PRODUCED (OR ARE ADMISSIBLE)

SEND PONG ONLY WHEN PING IS RECEIVED



$P_0, \text{SEND}(P_1, \text{PING})$



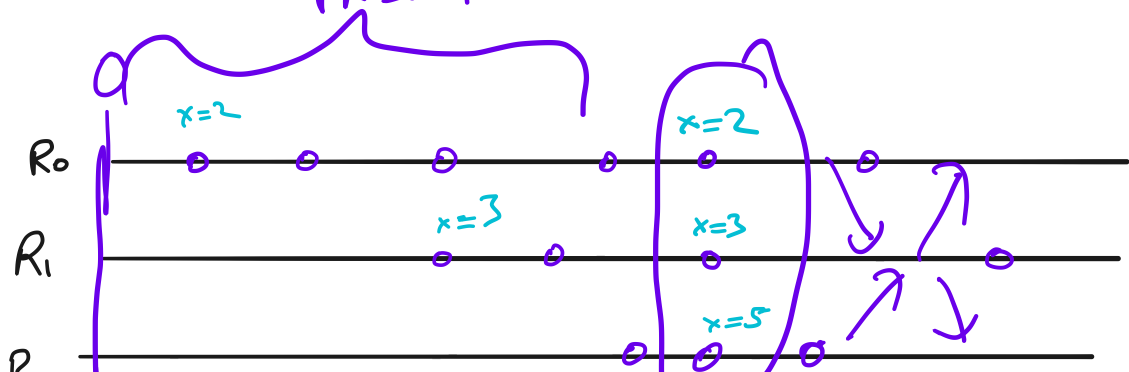
TYPES OF RESTRICTIONS / CORRECTNESS PROPERTIES

- SAFETY
↳ SOME BEHAVIOR (SEQ. OF EVENTS) NEVER OCCURS

LIVENESS

- SOME PROPERTY CAN/WILL EVENTUALLY HOLD

"EVENTUALLY ALL REPLICAS AGREE ON THE SET OF KEYS & VALUES"
PREFIX



$x=5$
 $\cdot \text{lock}(m);$
 $\text{if } (x \cdot 2 == 0) ?$
 $\quad \text{return}$
 else
 $\quad \text{unlock}(m);$
 $\quad \text{return}$

UNIFORM / ABSOLUTE ~~LIVENESS~~

↳ GIVEN A TRACE PREFIX, WHAT STEPS NEED TO BE TAKEN TO MEET PROPERTY

◦ GENERAL ADMISSIBLE
 FOR ANY PREFIX P , CAN FIND SUFFIX S_P
 S.T. CONDITION HOLDS FOR PS_P

◦ UNIFORM ADMISSIBLE
 THERE IS A SUFFIX S , S.T., FOR ANY PREFIX P
 CONDITION HOLDS FOR PS

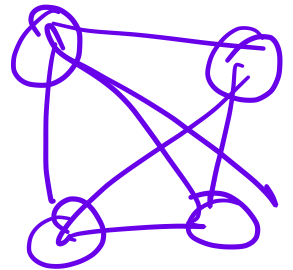
◦ ABSOLUTE
 GIVEN A TRACE T FOR WHICH CONDITION HOLDS &
 ANY ADMISSIBLE PREFIX P
 CONDITION HOLDS FOR PT

HARDER/EASIER?

CHANDY-LAMPORT

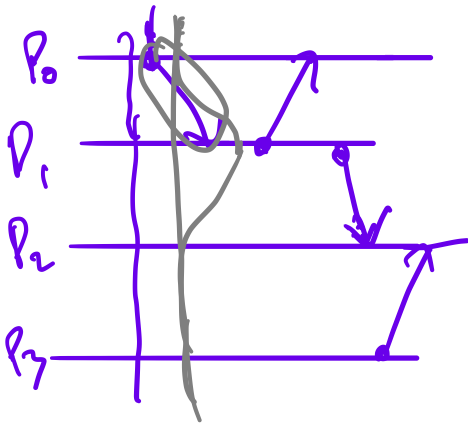
◦ GOAL & PROPERTIES?

◦ Compute ^{System's} global state



Snapshot

↳ Process
↳ Message channels



ASSUMPTIONS

- Ordered channels + reliable
- No failures

How?

- ① Save local state (1)
- ② Broadcast marker (1)
- ③ Process x on receiving

marker sends
back snapshot

