# The ABCDs of Paxos

Consensus: a set of processes *decide* on an input value

Main application: Replicated state machines

Paxos asynchronous consensus algorithm

    AP  Abstract Paxos:    generic, non-local version


    CP  Classic Paxos:    stopping failures, compare-and-swap
            1989: Lamport, Liskov and Oki

    DP  Disk Paxos:       stopping failures, read-write
            1999: Gafni and Lamport

    BP  Byzantine Paxos:  arbitrary failures
            1999: Castro and Liskov


*The paper and slides are at research.microsoft.com/lampson*

# Replicated State Machines

Lamport 1978: *Time, clocks and the ordering of events …*

Cast your problem as a deterministic state machine

    Takes client input requests for state transitions, called *steps*

    Performs the steps

    Returns the output to the client.

Make $n$ copies or 'replicas' of the state machine.

Use consensus to feed all the replicas the same inputs.

Steps must be deterministic, local to replica, atomic (use transactions)

Recover by replaying the steps (like transactions)

Even a read needs a step, unless the result is "as of step $n$".

# Applications of RSM

Reliable, available data storage system

Airplane flight control

**Reflexive applications:**

Changing quorums of the consensus algorithm

Issuing a *lease*:

    A lock on part of the state that times out, hence is fault tolerant

    Leaseholder can work on its state without consensus

    Like any lock, a lease can have modes or be hierarchical

# The Idea of Paxos

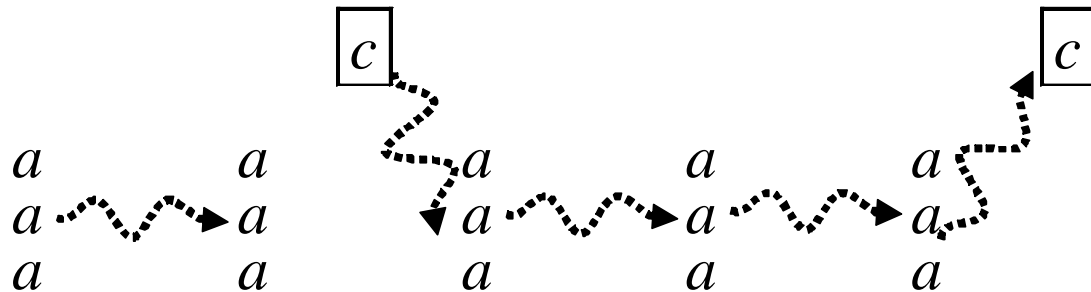A sequence of *views*; get a decision quorum in one of them.

Each view *v chooses* an *anchored* value $c_v$, equal to any earlier decision.

   If a quorum *accepts* the choice, decision!

Decision is irrevocable, may  be invisible, but is any later view's choice.
Choice   is changeable, must be visible if there was a decision

**Processes**

**Actions**

*Start*;
*Close*[a]   *Anchor*   *Choose*   *Accept*[a]   *Finish*[a];
                                                    *STEP*[a]

*Input*;

**Transmit**

$r^a$   **INPUT**   $c_v$   $r_v{}^a$   **OUTPUT**

**view change** | **normal operation**

# Design Methodology

- Communicate only *stable* predicates: once true always true

- Structure the program as a set of atomic actions

- Make actions as non-deterministic as possible: weakest guards

    Allows more freedom for the implementation
    Makes it clear what is essential

- Separate safety, liveness, and performance

    Safety first, then strengthen guards for liveness and scheduling

- Abstraction functions and simulation proofs

# Notation

Subscripts and superscripts for function arguments: $r_v{}^a$ for $r(v, a)$

State functions used like variables

Actions described like this:

| Name | Guard | State change |
|------|-------|--------------|
| **Close$_v$** | $c_v = nil \wedge x \in anchor_v$ | $\rightarrow c_v := x$ |

# Failure Model

A set $M$ of processes (machines)

A *faulty* process can send arbitrary messages: $F^m$

A *stopped* process does nothing: $S^m$

A *failed* process is faulty or stopped. State freezes after failure.

Limits on failure:

$Z_F$ = set of sets of processes that can all be faulty
$Z_S$ = set of sets of processes that can all be stopped
$Z_{FS}$ = set of sets of processes that can all be failed

Examples:

Fail-stop: $n$ processes,  $Z_F$={}, $Z_S$=$Z_{FS}$=any set of size < $(n+1)/2$

Byzantine: $n$ processes, $Z_F$ = $Z_S$=$Z_{FS}$=any set of size < $(n+1)/3$

Intel-Microsoft: $n_I + n_M$ processes, $Z_F$=any subset of one side

# Quorums and Predicates

Quorum set $Q$: set of sets of processes; $q$ in $\Rightarrow$ any superset in.

State predicate $g$. Predicate on processes $G$, so $G^m$ is a predicate.
    A *stable* predicate once true remains true.

$Q\#G$: A predicate $G$ appears to hold in quorum $Q$, $\{m \mid G^m \vee F^m\} \in Q$
    Shorthand: $Q[r_v{}^*{=}x]$ for $Q\#(\lambda\, m \mid r_v{}^m = x)$.

A *good* quorum is not all faulty: $Q_{\sim F} = \{q \mid q \notin Z_F\}$

$Q_1$ and $Q_2$ *exclusive*: $Q_1$ quorum for $G \Rightarrow$ no $Q_2$ quorum for its negation.

    Means $q_1 \cap q_2 \in Q_{\sim F}$ for any $q_1$ and $q_2$. Example: size $> (n + f)/2$

    Lift local $r_v{}^a{=}x \Rightarrow \sim(r_v{}^a{=}out)$ to global $Q_1[r_v{}^*{=}x] \Rightarrow \sim Q_2[r_v{}^*{=}out]$

$Q^+$: ensures $Q$ even after failures: $q^+ - z_{FS} \in Q$ for any $q^+, z_{FS}$

    A *live* quorum has $Q^+ \neq \{\}$

# Specification for Consensus

**type** $X$      $=...$      values to decide on

**var** $d$      $: (X \cup \{nil\}) := nil$    Decision

       $input$   $: \mathbf{set}\ X := \{\}$

| Name | Guard | State change |
|------|-------|--------------|
| *Input*($x$) | | $input := input \cup \{x\}$ |
| *Decision*: $X$ | $d \neq nil$ | $\rightarrow \mathbf{ret}\ d$ |
| | | |
| *Decide* | $d = nil \wedge x \in input \rightarrow d := x$ | |

# The Idea of Paxos

A sequence of *views*; get a decision quorum in one of them.
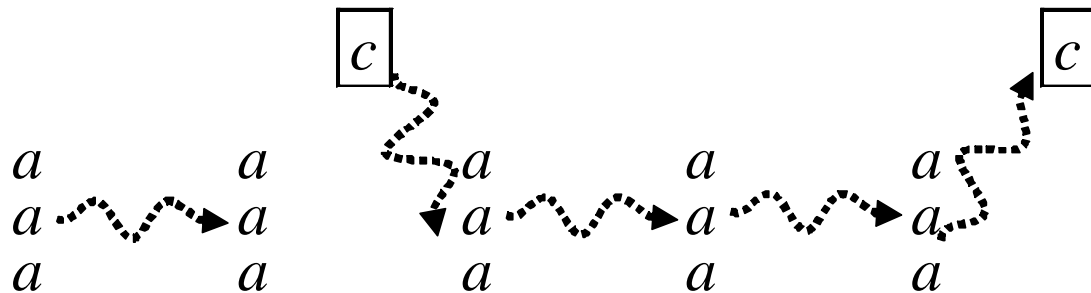
Each view *v chooses* an *anchored* value $c_v$: equals any earlier decision.

    If a quorum *accepts* the choice, decision!

<span style="color:green">Decision</span> is irrevocable, may  be invisible, but is any later view's <span style="color:red">choice</span>.
<span style="color:red">Choice</span>    is changeable, must be visible to *Anchor* if there was a <span style="color:green">decision</span>.



**Processes**

**Actions**

*Start*;
*Close*$^a$    *Anchor*    *Input*;
    <span style="color:red">*Choose*</span>    *Accept*$^a$    <span style="color:green">*Finish*$^a$</span>;
    **STEP**$^a$

**Transmit**    $r^a$    **INPUT**    <span style="color:red">$c_v$</span>    $r_v^{\,a}$    <span style="color:green">**OUTPUT**</span>

*view change*    *normal operation*

# Abstract Paxos—AP: State

State

| Non-local | Agents' | State functions | | | View is |
|---|---|---|---|---|---|
| | | | $r_v$ | $\underline{d}$ | |



$c_v$    1: $\left( \begin{array}{c} r_v^{\,1} \\ d^{\,1} \end{array} \right)$      $Q_{dec}[r_v^{\,*}=x]$      $x$    $x$     decided

*input*    2: $\left( \begin{array}{c} r_v^{\,2} \\ d^{\,2} \end{array} \right)$      $Q_{out}[r_v^{\,*}=out]$    *out*    *nil*     out

*active$_v$*    3: $\left( \begin{array}{c} r_v^{\,3} \\ d^{\,3} \end{array} \right)$

**else**          *nil*    *nil*     open
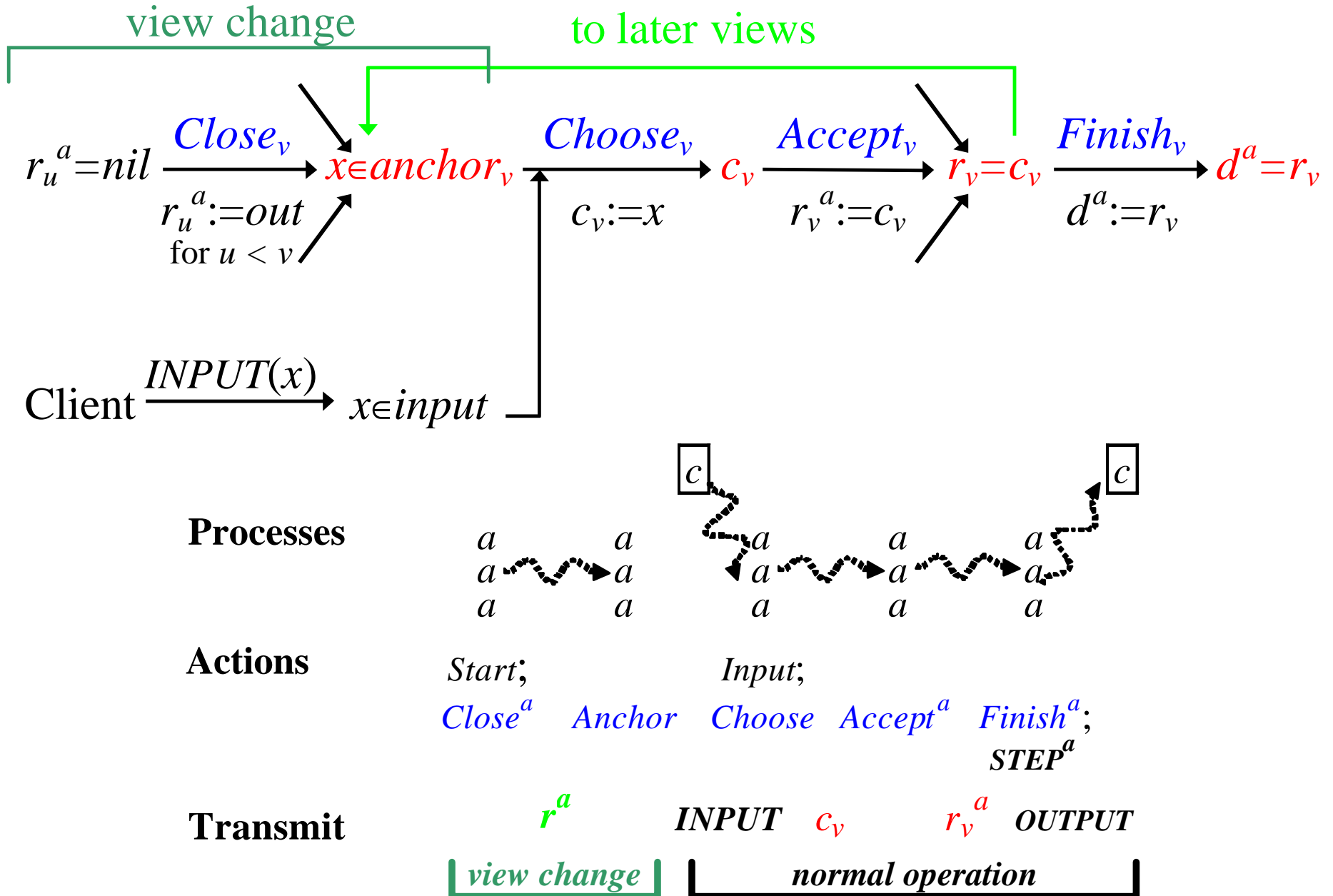
$Q_{dec}$ and $Q_{out}$ exclusive

*var = const* is stable for all these except *input*, and $x \in input$ is stable.

# AP: Data Flow

view change        to later views

$$r_u^a = nil \xrightarrow[\substack{r_u^a := out \\ \text{for } u < v}]{Close_v} x \in anchor_v \xrightarrow[c_v := x]{Choose_v} c_v \xrightarrow[r_v^a := c_v]{Accept_v} r_v = c_v \xrightarrow[d^a := r_v]{Finish_v} d^a = r_v$$

$$\text{Client} \xrightarrow{INPUT(x)} x \in input$$

**Processes**

$$\begin{matrix} a & & a & & c & & a & & a & & c \\ a & & a & & a & & a & & a & & a \\ a & & a & & a & & a & & a & & a \\ a & & a & & a & & a & & a & & a \end{matrix}$$

**Actions**

$Start$;        $Input$;

$Close^a$    $Anchor$    $Choose$   $Accept^a$    $Finish^a$;

$STEP^a$

**Transmit**

$r^a$        $INPUT$   $c_v$        $r_v^a$   $OUTPUT$

view change        normal operation

# Example

| | $c_v$ | $r_v^a$ | $r_v^b$ | $r_v^c$ | $c_v$ | $r_v^a$ | $r_v^b$ | $r_v^c$ |
|---|---|---|---|---|---|---|---|---|
| View 1 | 7 | 7 | *out* | *out* | 8 | 8 | *out* | *out* |
| View 2 | 8 | *out* | 8 | *out* | 9 | 9 | *out* | 9 |
| View 3 | 9 | *out* | *out* | 9 | 9 | *out* | *out* | 9 |
| *input* ∩ *anchor*$_4$ | = {7, 8, 9} seeing *a, b, c* <br> ⊇{8}　　seeing *a, b* <br> ⊇{9}　　seeing *a, c* or *b, c* | | | | {9} no matter what <br> quorum we see | | | |

Two runs of AP with
   agents *a, b, c*,
   two agents in a quorum,
   *input* = {7, 8, 9}

# Anchoring

**invariant** $r_v = x \wedge r_u = x' \Rightarrow x = x'$      all results agree

$= \ \forall \, x', u \mid r_v = x \wedge \underline{r_u = x'} \Rightarrow x = x'$      assume $u < v$

$= \ r_v = x \Rightarrow (\forall \, u < v, x' \neq x \mid \sim \boxed{Q_{dec}[r_u{}^* = x']})$      ${\color{red} r_u{}^a \in \{x,\, out\}}$

                                                ${\color{red} \Rightarrow \sim(r_u{}^a = x')}$

$\Leftarrow \ r_v = x \Rightarrow (\forall \, u < v \mid {\color{red} Q_{out}[r_u{}^* \in \{x, out\}]})$

**sfunc** $anchor_v$

$= \quad\quad \{x \mid (\forall \, u < v \mid \quad\quad\quad\quad Q_{out}[r_u{}^* \in \{x, out\}])\}$

$= \quad\quad \{x \mid (\forall \, w \mid v_0 = w < u \Rightarrow Q_{out}[r_w{}^* \in \{x, out\}])\}$      $= anchor_u$

$\quad\quad \cap \{x \mid \quad\quad\quad\quad\quad\quad Q_{out}[r_u{}^* \in \{x, out\}]\}$

$\quad\quad \cap \{x \mid (\forall \, w \mid u \ < w < v \Rightarrow Q_{out}[r_w{}^* \in \{x, out\}])\}$      $= X$ **if** $out_{u,v}$

$= \ anchor_u \cap \{x \mid Q_{out}[r_u{}^* \in \{x, out\}]\}) \ \textbf{if} \ out_{u,v}$

$\supseteq \ \textbf{if} \ out_{v_0,v} \ \textbf{then} \ X \ \textbf{elseif} \ out_{u,v} \wedge r_u{}^a = x \ \textbf{then} \ \{x\} \ \textbf{else} \ \{\,\}$

    where $out_{u,v} = (\forall \, w \mid u < w < v \Rightarrow r_w = out)$

# AP: Algorithm

$Start_v$ $\quad u<v$ too slow $\quad\rightarrow active_v := true$

$Close_v^a$ $\quad active_v$ $\qquad\rightarrow$**for all** $u < v$ **do** $\qquad$ **post** $u<v$
$\qquad\qquad\qquad\qquad\qquad\quad$ **if** $r_u^a = nil$ $\qquad\qquad \Rightarrow r_u^a \neq nil$
$\qquad\qquad\qquad\qquad\qquad\quad$ **then** $r_u^a := out$

$anchor_v = anchor_u \cap \{x \mid Q_{out}[r_u^* \in \{x,out\}]\})$ **if** $out_{u,v}$

$Anchor_v$ $\quad anchor_v \neq \{\}$ $\qquad\rightarrow$no state change

$Choose_v$ $\quad c_v^a = nil$ $\qquad\qquad\rightarrow c_v := x$
$\qquad\qquad \wedge x \in input \cap anchor_v$

$Accept_v^a$ $\quad r_v^a = nil$ $\qquad\qquad\rightarrow r_v^a := c_v;\ Close_v^a$
$\qquad\qquad \wedge c_v \neq nil$

$Finish_v^a$ $\quad r_v \in X$ $\qquad\qquad\rightarrow d^a := r_v$

to later views

$r_u^a=nil \xrightarrow{\ Close_v\ } x \in anchor_v \xrightarrow{\ Choose_v\ } c_v \xrightarrow{\ Accept_v\ } r_v=c_v \xrightarrow{\ Finish_v\ } d^a=r_v$
$\qquad\qquad r_u^a:=out$ $\qquad\qquad\quad c_v:=x \qquad\qquad\quad r_v^a:=c_v \qquad\qquad d^a:=r_v$
$\qquad\qquad$ for $u < v$

15

# AP: Liveness

*Choose$_v$* must see an element of *input* $\cap$ *anchor$_v$*.

Recall *anchor$_v$*

$= \quad anchor_u \cap \{x \mid Q_{out}[r_u^* \in \{x, out\}]\}$          **if** *out$_{u,v}$*

$\supseteq$   **if** *out$_{v0,v}$* **then** *X* **elseif** *out$_{u,v}$* $\wedge\ r_u^a = x$ **then** $\{x\}$ **else** $\{\ \}$

After *Close$_v^a$*, an OK agent *a* has $r_u^a \neq nil$ for all $u < v$.

So if *Q$_{out}$* is live, we see either $u < v$ is out, or $r_u^a = x$ for some OK *a*.

But $r_u^a = c_u \in$ *input* $\cap$ *anchor$_u$*

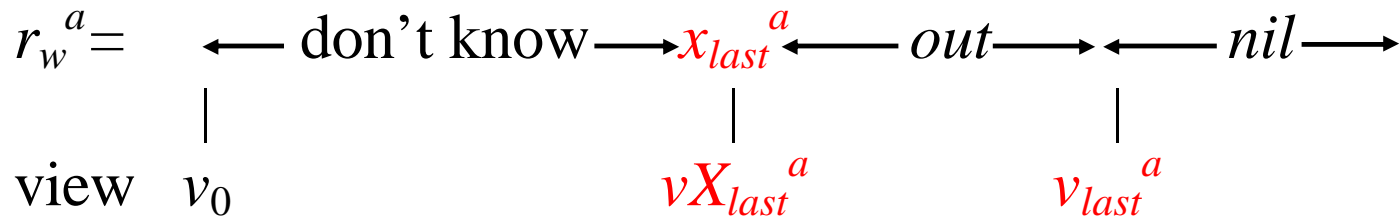     If we know *a* is OK, then $r_u^a$ is what we want

     With faults (in BP), we might not know.
     But if *anchor$_u$* is visible, that is enough.

Still not live if new views start too fast.

# Optimizations

**Fixed-size agent state:**

$$r_w{}^a = \quad \longleftarrow \text{don't know} \longrightarrow x_{last}{}^a \longleftarrow out \longrightarrow \longleftarrow nil \longrightarrow$$

view $\quad v_0 \qquad\qquad\qquad vX_{last}{}^a \qquad\qquad v_{last}{}^a$

**Successive steps:**

Because $anchor_v$ doesn't depend on *input*, can compute it for lots of steps at once.

This is called a *view change*

One view change is enough for any number of steps

Can batch steps, with one Paxos/batch.

Can run steps in parallel, subject to external consistency.

# Disk Paxos—DP

The goal—Replace the conditional writes in *Close* and *Accept* with simple writes.

$\textbf{\textit{Accept}}_v^a$ $\quad r_v^a = nil \land c_v \neq nil \quad \to r_v^a := c_v;\ \textit{Close}_v^a$

The idea—Replace $r_v^a$ with $rx_v^a$ and $ro_v^a$.

$\textbf{\textit{Accept}}_v^a$ $\quad c_v \neq nil \quad\quad\quad\quad\quad\quad \to rx_v^a := c_v;\ \textit{Close}_v^a$

$\textbf{\textit{Close}}_v^a$ $\quad active_v \quad\quad\quad\quad\quad\quad \to \textbf{for all } u < v \textbf{ do } ro_u^a := out$

Proof: Keep $r_v^a$ as a history variable. Abstract it to AP's $\underline{r}_v^a$.
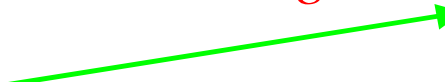This invariant makes it work (sometimes with an extra view).

| $rx_v^a =$ | $\land$ | $ro_v^a =$ | $\Rightarrow$ | $r_v^a$ |
|---|---|---|---|---|
| *nil* | | *nil* | | $= nil$ |
| *nil* | | *out* | | $= out$ |
| $x$ | | *nil* | | $= x$ |
| $x$ | | *out* | | $\neq nil$ |

# Communication

A process has knowledge $T$ of stable non-local facts

$$g@m = (T^m \Rightarrow g)$$

We transmit these facts (note that transmitter $k$ may be failed):

$\boldsymbol{Transmit^{k,m}(g)} \quad g@k \lor F^k \rightarrow T^m := T^m \land (g@k \lor F^k) \ \textbf{post} \ (g@k \lor F^k)@m$

A faulty $k$ can transmit anything:

A fact known to a $Q_{\sim F}^+$ quorum is henceforth known to a $Q_{\sim F}$ quorum of OK agents, and therefore eventually known to everyone.

$\boldsymbol{Broadcast^m(g)} \quad Q_{\sim F}^+[g@*] \land OK^m \rightarrow T^m := T^m \land g \qquad \textbf{post} \ g@m$

Implement $Transmit^{k,m}$ by sending messages. It's fair if $k$ is OK. This works because the facts are stable.

# Classic Paxos—CP

The goal—Tolerate stopped processes

The idea—Agents are the same as in AP. Use a *primary* process to:
    Implement *Choose*
    Compute an estimate $re_v$ of $r_v$
    Relay facts among the agents
    Do all the scheduling.

So the primary sends $active_v$ to agents to enable $Close_v$, collects $r^a$, computes *anchor*, gets inputs, does *Choose*, sends $c^p$ to agents, collects $r^a$ again to compute $re_v$, and sends $d$.

**Choose$^p$**        $active^p \wedge c^p = nil$        $\rightarrow c^p := x$
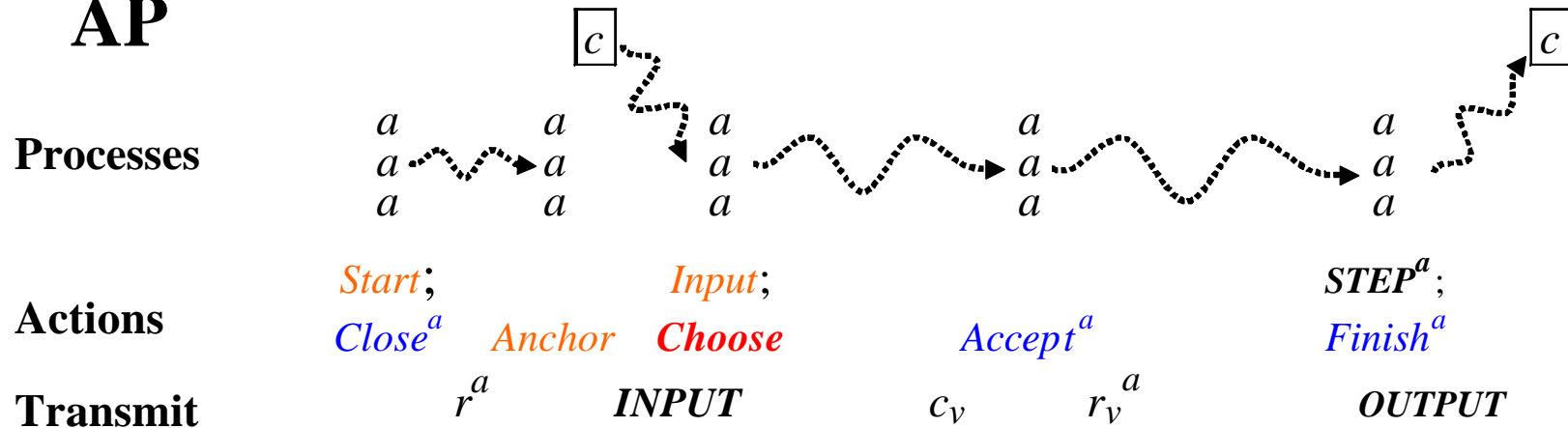            $\wedge\, x \in input^p \cap anchor^p$

Must have only one $c^p$ per view. Get this with
    At most one primary per view, and
    Primary chooses at most once per view

# AP and CP

## AP



**Processes**

**Actions**

Start; Close$^a$   Anchor   **Choose**   Input;   Accept$^a$   **STEP$^a$**; Finish$^a$

**Transmit**   $r^a$   **INPUT**   $c_v$   $r_v^a$   **OUTPUT**

## CP



**Processes**

**Actions**   Start$^p$; Close$^p$ Close$^a$   Anchor$^p$ **Choose$^p$**; Accept$^p$ Accept$^a$   Input$^p$;   **STEP$^p$** Finish$^p$;   **STEP$^a$** Finish$^a$;

**Transmit**   $active_v$   $r^a$   **INPUT**   $c^p$   $r_v^a$   **OUTPUT**   $re_v^p$

**Messages**   $1 \rightarrow n^*$   $n^* \rightarrow 1$   $\boxed{1 \rightarrow 1}$   $1 \rightarrow n^*$   $n^* \rightarrow 1$   $\boxed{1 \rightarrow 1}$   $1 \rightarrow n^*$

**Primary:**   Relay   Choose $c_v$   Estimate $r_v$

# Byzantine Paxos—BP

The goal—Tolerate faulty processes

The idea—To ensure one $c_v$, a self-exclusive <span style="color:red">quorum $Q_{ch}$</span> chooses it

Still have a primary to <span style="color:blue">propose $c_v$</span>; an OK agent only chooses this

Primary's proposal should be anchored and input at agents
A faulty primary can stop its view from deciding

Every agent needs an estimate $ce_v^a$ of $c_v$ and an estimate $re_v^a$ of $r_v$

Invariant: The estimates either are *nil* or equal the true values.

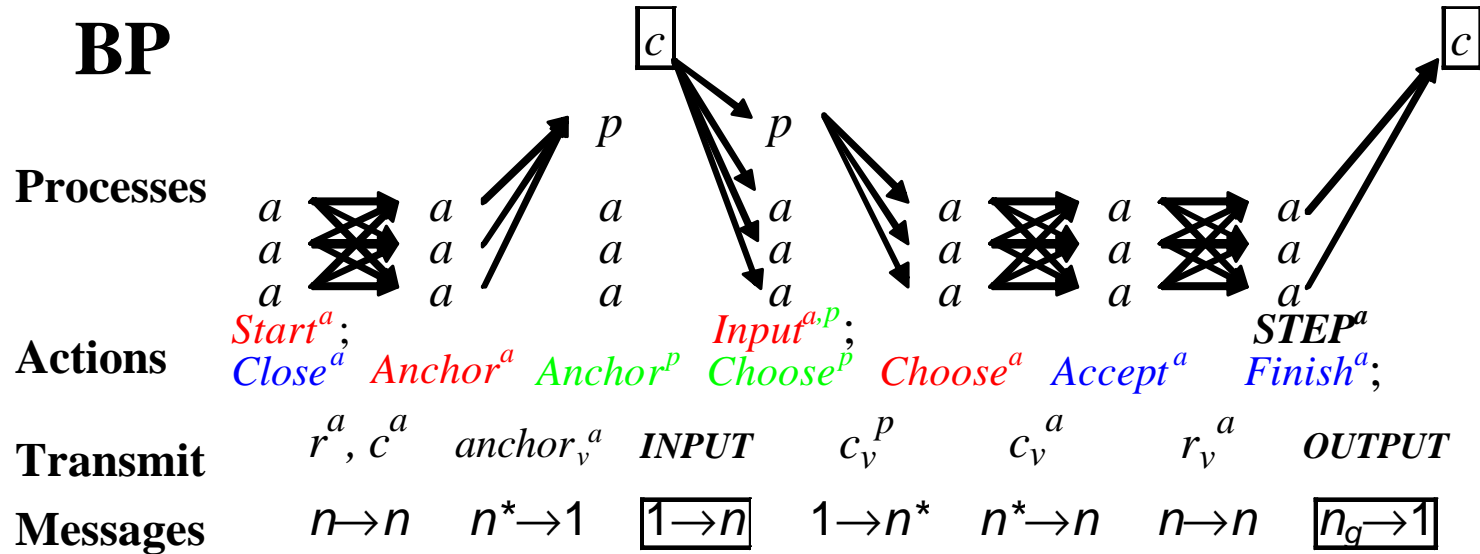Every agent also needs its own *input*[a]

**abstract** $\underline{c_v}$ $= \textbf{if}$ $\color{red}{Q_{ch}[c_v{}^*{=}x]}$ $\textbf{then } x$ $\textbf{else } nil$

**sfunc** $ce_v^a = \textbf{if}$ $(Q_{ch}[c_v{}^*{=}x])@a$ $\textbf{then } x$ $\textbf{else } nil$

$anchor_v^a = anchor_u \cap \{x \mid Q_{out}[r_u{}^* {\in} \{x, out\}]@a\}$ $\textbf{if } out_{u,v}^a$

$\color{blue}{anchor_v^p = \{x \mid Q_{\sim F}{}^+[x {\in} anchor_v{}^*]@p\}}$

# CP and BP

## CP

**Processes**

*p*   *p*   $c$   *p*   *p*   *p*   $c$

*a a a*   *a a a*   *a a a*   *a a a*   *a a a*   *a a a*

**Actions**

$STEP^p$   $STEP^a$

$Start^p$ ;   $Input^p$ ;

$Close^p$ $Close^a$   $Anchor^p$ $Choose^p$ ; $Accept^p$ $Accept^a$   $Finish^p$ ;   $Finish^a$ ;

**Transmit**   $active_v$   $r^a$   **INPUT**   $c^p$   $r_v^a$   **OUTPUT**   $re_v^p$

**Messages**   $1{\to}n*$   $n*{\to}1$   $\boxed{1{\to}1}$   $1{\to}n*$   $n*{\to}1$   $\boxed{1{\to}1}$   $1{\to}n*$

## BP

**Processes**

$c$   *p*   *p*   $c$

*a a a*   *a a a*   *a a a*   *a a a*   *a a a*   *a a a*   *a a a*

**Actions**

$STEP^a$

$Start^a$ ;   $Input^{a,p}$ ;

$Close^a$   $Anchor^a$   $Anchor^p$   $Choose^p$   $Choose^a$   $Accept^a$   $Finish^a$ ;

**Transmit**   $r^a, c^a$   $anchor_v^a$   **INPUT**   $c_v^p$   $c_v^a$   $r_v^a$   **OUTPUT**

**Messages**   $n{\to}n$   $n*{\to}1$   $\boxed{1{\to}n}$   $1{\to}n*$   $n*{\to}n$   $n{\to}n$   $\boxed{n_q{\to}1}$

23

# Liveness of BP

*Choose* must see an element of *input* $\cap$ *anchor$_v$*.

Recall $anchor_v \supseteq anchor_u \cap \{x \mid Q_{out}[r_u^* \in \{x, out\}]\}$

After $Close_v^a$, an OK agent $a$ has $r_u^a \neq nil$ for all $u < v$.

So if $Q_{out}$ is live, we see either $u < v$ is out, or $r_u^a = x$ for some OK $a$.

But $r_u^a = c_u \in input \cap anchor_u$

    Unfortunately, we don't know whether $a$ is OK.

    But we do have $Q_{ch}[c_u^* = x]$, hence $Q_{ch}[(x \in anchor_u)@a]$

    So if $Q_{ch}$ is live, $x \in anchor_u$ is broadcast, which is enough.

So either we eventually see all previous views out,
or we see $x \in anchor_u$ and all views between $u$ and $v$ out.

A faulty client can wreck a view by not sending input to all agents.

# Conclusion

Paxos is a practical protocol for fault-tolerant asynchronous consensus.

Paxos is efficient in replicated state machines, which are the best mechanism for most fault-tolerant systems.

Paxos works in a sequence of views,

    Each view chooses a value and then seeks a decision quorum.

    A later view chooses any possible earlier decision

Abstract Paxos chooses a consensus value non-locally, and then decides by local actions of the agents.

    The agents are read-modify-write memories.

    Disk Paxos generalizes this to read-write memories.

Classic Paxos uses a primary process to choose.

Byzantine Paxos uses a primary to propose, a quorum to choose.