# Mutually Assured Destruction[†] (or the Joy of Sync)

Peter Welch (`phw@kent.ac.uk`),
Jan Bækgaard Pedersen (`matt.pedersen@unlv.edu`)
Frederick R.M. Barnes (`frmb@kent.ac.uk`)

CPA 2013 Fringe, Napier University, 25 August, 2013

[†] plus non-blocking barriers and performance …
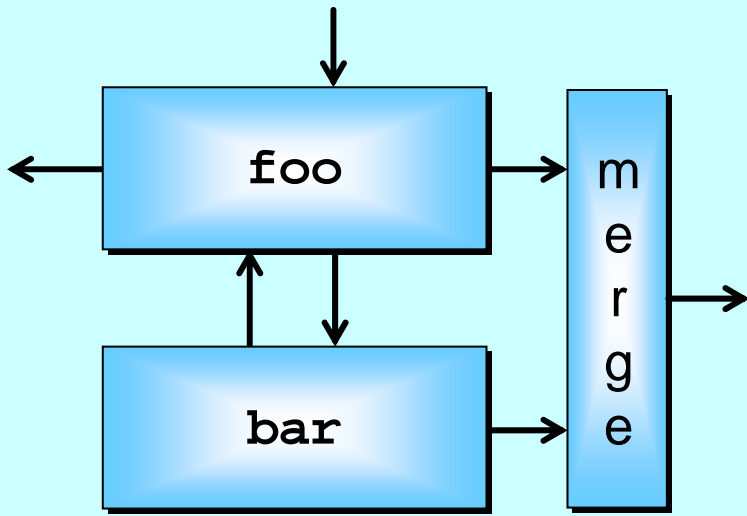
# The Joy of Sync

Process oriented design ...

Synchronous communications …

Synchronous barriers …
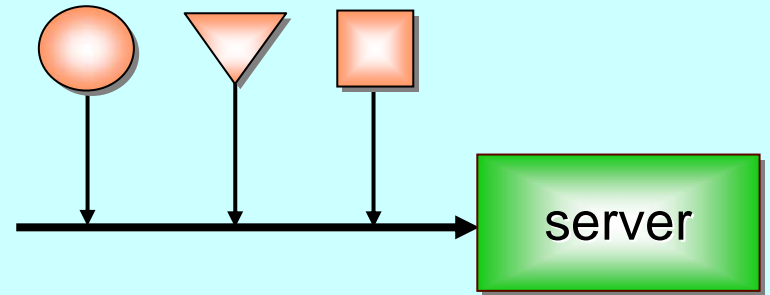
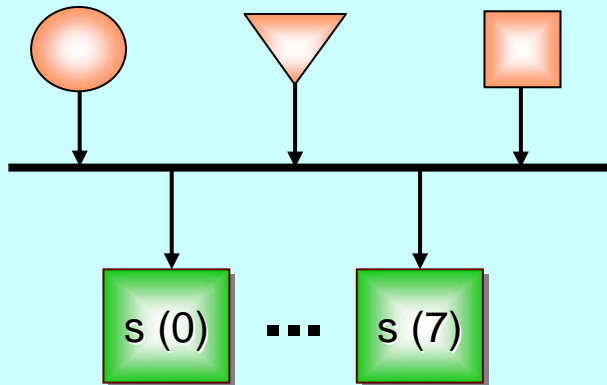Mutually assured destruction …

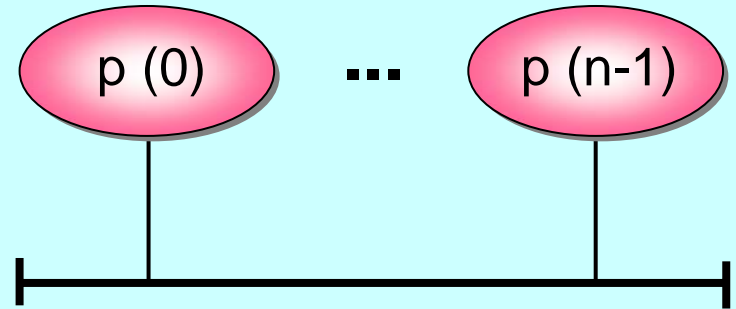Non-blocking barriers …

Performance …

**(a) a network of three processes, connected by four internal (hidden) and three external channels.**
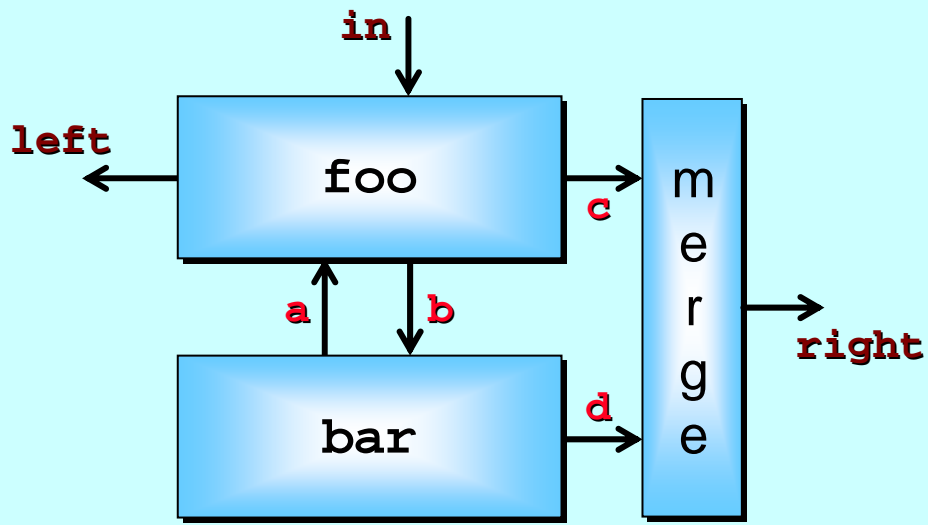
**(b) three processes sharing the writing end of a channel to a server process.**

**(c) three processes sharing the writing end of a channel to a bank of servers sharing the reading end.**

**(d) n processes enrolled on a shared barrier (any process synchronising must wait for all to synchronise).**

**(a) a network of three processes, connected by four internal (hidden) and three external channels.**

```
CHAN BYTE a, b, c, d:
PAR
  foo (in?, left!, a?, b!, c!)
  bar (a!, b?, d!)
  merge (c?, d?, right!)
```

```
PROC thing (CHAN INT in?, left!, right!)
  CHAN BYTE a, b, c, d:
  PAR
    foo (in?, left!, a?, b!, c!)
    bar (a!, b?, d!)
    merge (c?, d?, right!)
:
```

**in**

**left**

**right**

**thing**

process
abstraction

```
PROC thing (CHAN INT in?, left!, right!)
  CHAN BYTE a, b, c, d:
  PAR
    foo (in?, left!, a?, b!, c!)
    bar (a!, b?, d!)
    merge (c?, d?, right!)
:
```
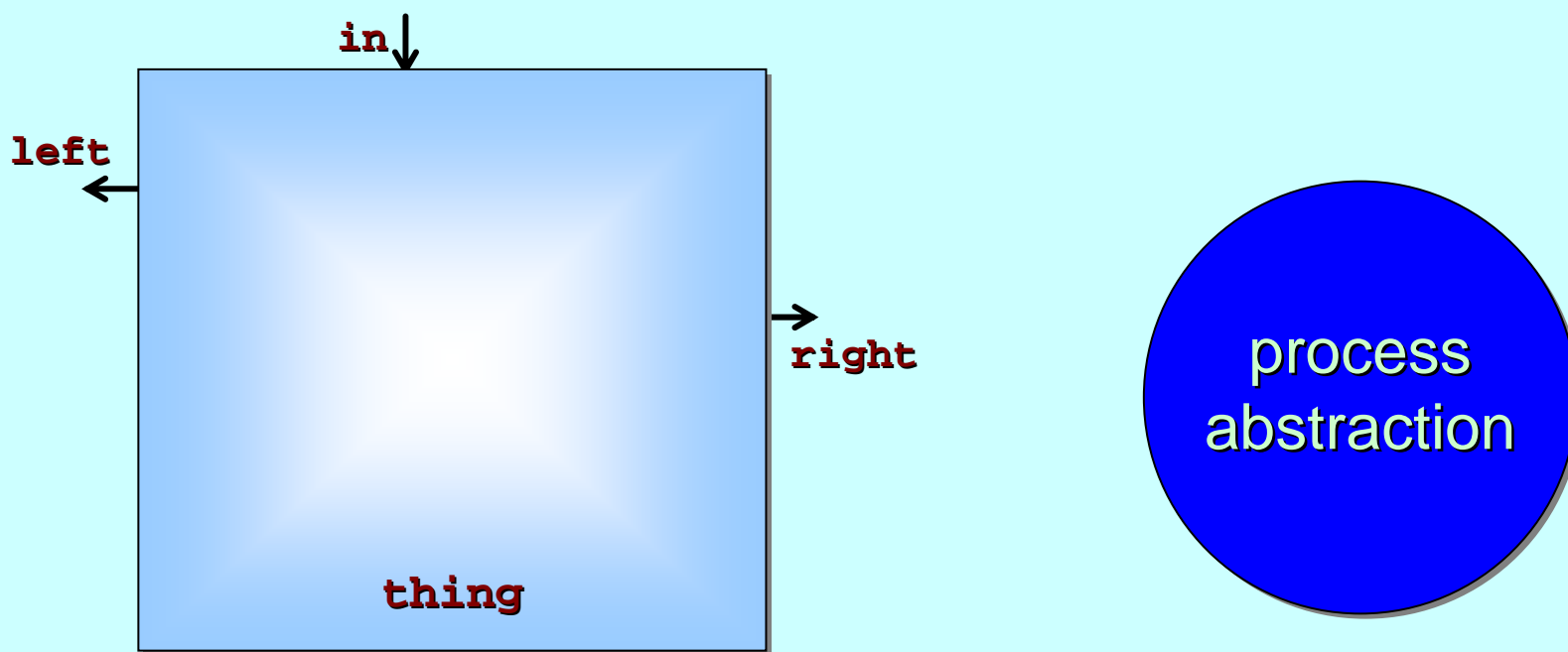
**in**

**left**  **thing**  **right**

process
abstraction

```
PROC thing (CHAN INT in?, left!, right!)
```

Like **foo**, **bar** and **merge** previously, **thing** is a process that can be used as a component in another network.

Concurrent systems have structure – networks within networks. We must be able to express this!  And we can … ☺ ☺ ☺

**(b)** **three processes sharing the writing end**
**of a channel to a server process.**

```
SHARED ! CHAN SOME.SERVICE c:
PAR
  circle (c!)
  triangle (c!)
  square (c!)
  server (c?)
```

```
SHARED CHAN ANOTHER.SERVICE c:
PAR
  PAR
    circle (c!)
    triangle (c!)
    square (c!)
  PAR i = 0 FOR 8
    s (i, c?)
```



c

s (0) ▪▪▪ s (7)

**(c) three processes sharing the writing end of a channel
  to a bank of servers sharing the reading end.**

```
BARRIER b:
PAR i = 0 FOR n ENROLL b
  p (i, b)
```



(d) n processes enrolled on a shared barrier (any process synchronising must wait for all to synchronise).

# The Joy of Sync

Process oriented design ...

Synchronous communications …

Synchronous barriers …

Mutually assured destruction …

Non-blocking barriers …

Performance …

# Synchronised Communication

A

B

c

c ! 42

c ? x

**A** may *write* on **c** at any time, but has to wait for a *read*.

**B** may *read* from **c** at any time, but has to wait for a *write*.

( A (c) ‖ B (c) )  \  {c}

# Synchronised Communication

A ⟶ c ⟶ B

A
c ! 42

B
c ? x

Only when both **A** and **B** are ready can the communication proceed over the channel **c**.

( A (c) ‖ B (c) )  \  {c}

# Synchronised Communication

$$A \xrightarrow{\quad c \quad} B$$

- **Benefit**
  - ◆ Once the writer has written, it *knows* the reader has read

OK: plenty of other processes to run and ultra-fast context switch (comparable to a procedure call)

- **Careful**
  - ◆ Writer blocks if reader is not ready
  - ◆ Lots of deadlock possibilities

OK: work with (a small set of) synchronisation patterns for which we have proven safety theorems

# Simple Deadlock Example



If there is no discipline on when *A* and *B* communicate, then *A* may commit to output on *c*, followed by *B* on *d* … or vice-versa.  Either way, neither are listening and both are stuck.  Same happens if both commit to input.

☹ ☹ ☹ ☹ ☹

# Client-Server Pattern



*client*: makes a *request* any time, then commits to taking *reply*.

*server*: always accepts a *request* (within some bounded time), then always makes a *reply* (within some bounded time).  It may make requests itself, as a *client* to other *servers*.

No deadlock is now possible from this client-server relationship.

☺ ☺ ☺ ☺ ☺

# Client-Server Pattern



*client*: makes a *request* any time, then commits to taking *reply*.

*server*: always accepts a *request* (within some bounded time), then always makes a *reply* (within some bounded time). It may make requests itself, as a *client* to other *servers*.

**Symbology:** this represents a client-server relation. It points *to* the server and allows a *2-way* conversation (initiated by the client)

# Client-Server Pattern

A *server* may have many *clients* …



Only one *client* at a time converses with the *server*.  They form an orderly queue.  Still no deadlock possible – and no *client* starvation.  No polling on the queue, so no livelock either.

# Client-Server Theorem

A *client-server* system that has no cycles in its *client-server* relations is deadlock, livelock and starvation free.

# The Joy of Sync

Process oriented design ...

Synchronous communications …

Synchronous barriers …

Mutually assured destruction …

Non-blocking barriers …

Performance …

# Barriers

The occam-π `BARRIER` type corresponds to a multiway CSP *event*, though some higher level design patterns (such as *resignation*) have been built in.



Basic CSP semantics apply. When a process *synchronises* on a barrier, it blocks until all other processes *enrolled* on the barrier have also *synchronised*. Once the barrier has completed (i.e. all *enrolled* processes have *synchronised*), all blocked processes are rescheduled for execution.

# Barriers

The occam-π **BARRIER** type corresponds to a multiway CSP *event*, though some higher level design patterns (such as *resignation*) have been built in.

| worker (0) | worker (1) | ... | worker (n-1) |
|---|---|---|---|

b

```
BARRIER b:
PAR i = 0 FOR n ENROLL b
  worker (i, b)
```

A **PAR** construct must *explicitly* **ENROLL** its components on barriers

The number of processes enrolled on an in-scope barrier is unchanged by a *non-enrolling* **PAR** – but *only one* of its components may reference it.

# Barriers

Processes may synchronise on more than one barrier:



```
BARRIER b, c:
PAR i = 0 FOR n ENROLL b, c
  worker (i, b, c)
```

To synchronise on a barrier:

```
SYNC b
```
or
```
SYNC c
```

# Barriers

Barriers are commonly used to synchronise multiple *phases* of computation between a set of processes.  Within each phase, other synchronisations (channel / barrier) may take place:

```
PROC worker (VAL INT id, BARRIER b, c)
   ...   local declarations / initialisation
   WHILE running
     SEQ
       SYNC b
       ...   observe neighbourhood phase
       SYNC c
       ...   change neighbourhood phase
:
```

**All workers do this together – all see the same thing …**

**All workers do this together – may need to negotiate …**

# **Barriers**

*But it's safer programming for each phase to be synchronised by its own barrier …*

Of course, only *one* barrier is actually needed to synchronise the phases in this example:

```
PROC worker (VAL INT id, BARRIER a)
  ...   local declarations / initialisation
  WHILE running
    SEQ
      SYNC a
      ...   observe neighbourhood phase
      SYNC a
      ...   change neighbourhood phase
:
```

**All workers do this together – all see the same thing …**

**All workers do this together – may need to negotiate …**

# Barriers – Safety

occam-π **BARRIER** *synchronisation* is *safe* in the sense that *enrollment* and *resignation* are automatically managed.  A process may *synchronise* on a **BARRIER** if and only if it is *enrolled*.

Try to break this rule … your program won't compile.  There are zero memory and run-time costs to enforce it.  ☺

# The Joy of Sync

Process oriented design ...

Synchronous communications …

Synchronous barriers …

Mutually assured destruction …

Non-blocking barriers …

Performance …

# Mutually Assured Destruction

Two processes are given, at the same time, their own task to complete; we are satisfied with the completion of either one of them; whichever process finishes first interrupts the other and reports its completion; the one that is interrupted abandons its task and reports that fact.

Such requirements are common in control systems, robotics, e-commerce, model-checking, …

**e.g.**
- Drive rover vehicle forwards *target* meters.
- Look out for *Martians*.
- Stop and report when either is achieved.

**motorSensor**

**cameraSensor**

```
monitor
(move)
```

**a**

**b**

```
monitor
(search)
```

**MADsystem**

**moveCommand**

**searchCommand**

**moveReport**

**searchReport**

**e.g.**

– Drive rover vehicle forwards *target* meters.

– Look out for *Martians*.

– Stop and report when either is achieved.

```
INT
```

**sensor**

```
monitor
(mode)
```

**killMe**

**killYou**

```
PROTOCOL KILL
  CASE
    kill
:
```

```
INT
```

**command**

**report**

```
PROTOCOL REPORT
  CASE
    me          -- task completed
    she         -- task abandoned
:
```

```
PROC monitor (VAL INT mode, CHAN INT command?, sensor?,
              CHAN REPORT report!,
              CHAN KILL killYou!, killMe?)
  WHILE TRUE
    PRI ALT
      INT target:
      command ? target                 -- service requested
        service (mode, target, sensor?, report!,
                 killYou!, killMe?)
      INT x:
      sensor ? x                        -- accept and discard
        SKIP
:
```

Copyright  P.H.Welch, (2013)

```
PROC service (VAL INT mode, target, CHAN INT sensor?,
              CHAN REPORT report!,
              CHAN KILL killYou!, killMe?)
  ...  local state and initialisation
  INITIAL BOOL running IS TRUE:
  WHILE running
    PRI ALT
      killMe ? kill
        ...  report 'she' and exit loop
      INT x:
      sensor ? x
        ...  process x
:
```

**sensor**

**monitor**
**(mode)**

**killMe**

**killYou**

**report**  **command**

```
killMe ? kill
  SEQ
    report ! she
    running := FALSE
```

```
PROC service (VAL INT mode, target, CHAN INT sensor?,
              CHAN REPORT report!,
              CHAN KILL killYou!, killMe?)
  ...  local state and initialisation
  INITIAL BOOL running IS TRUE:
  WHILE running
    PRI ALT
      killMe ? kill
        ...  report 'she' and exit loop
      INT x:
      sensor ? x
        ...  process x
:
```

```
    INT x:
    sensor ? x
      SEQ
        ...  update local state with x (depends on mode)
        IF
          ...  task complete
            SEQ
              killYou ! kill
              report ! me
              running := FALSE
          TRUE
            SKIP
  :
```

Copyright  P.H.Welch, (2013)

motorSensor               cameraSensor

monitor
(move)

**a**

**b**

monitor
(search)

**MADsystem**

moveCommand       searchCommand

moveReport         searchReport

```
PROC MADsystem (CHAN INT moveCommand?, searchCommand?,
                CHAN INT motorSensor?, cameraSensor?,
                CHAN REPORT moveReport!, searchReport!)
  CHAN KILL a, b:
  PAR
    monitor (move, moveCommand?, motorSensor?,
             moveReport!, b!, a?)
    monitor (search, searchCommand?, cameraSensor?,
             searchReport!, a!, b?)
:
```

**In Service**

average sensor data interval = 10 ms (varying)
average sensor inputs per service = 100 (varying)

motorSensor

cameraSensor

**MADsystem**

moveCommand

moveReport

searchCommand

searchReport

**Ran for 2 years (approx. 64m trials):   D E A D L O C K E D**   ☹ ☹ ☹

# Should have asked for a model check …

```
VERIFY DEADLOCK.FREE MADsystem
```
✗

motorSensor → **MADsystem** ← cameraSensor

moveCommand
moveReport
searchReport
searchCommand

**A trace leading to deadlock is provided:**

```
<moveCommand, motorSensor, searchCommand, cameraSensor>
```

# Should have asked for a model check …

A trace leading to deadlock is provided:

`<moveCommand, motorSensor, searchCommand, cameraSensor>`

motorSensor

cameraSensor

monitor
(move)

a

b

monitor
(search)

**MADsystem**

moveCommand

searchCommand

moveReport

searchReport

```occam
PROC monitor (VAL INT mode, CHAN INT command?, sensor?,
              CHAN REPORT report!,
              CHAN KILL killYou!, killMe?)
  WHILE TRUE
    PRI ALT
      INT target:
      command ? target              -- service requested
        service (mode, target, sensor?, report!,
                 killYou!, killMe?)
      INT x:
      sensor ? x                    -- accept and discard
        SKIP
:
```

Copyright  P.H.Welch, (2013)

```
sensor

        monitor          killMe
        (mode)          killYou

report    command
```

```
<moveCommand, motorSensor,
searchCommand, cameraSensor>
```

```occam
PROC service (VAL INT mode, target, CHAN INT sensor?,
              CHAN REPORT report!,
              CHAN KILL killYou!, killMe?)
  ...  local state and initialisation
  INITIAL BOOL running IS TRUE:
  WHILE running
    PRI ALT
      killMe ? kill
        ...  report 'she' and exit loop
      INT x:
      sensor ? x
        ...  process x
:
```
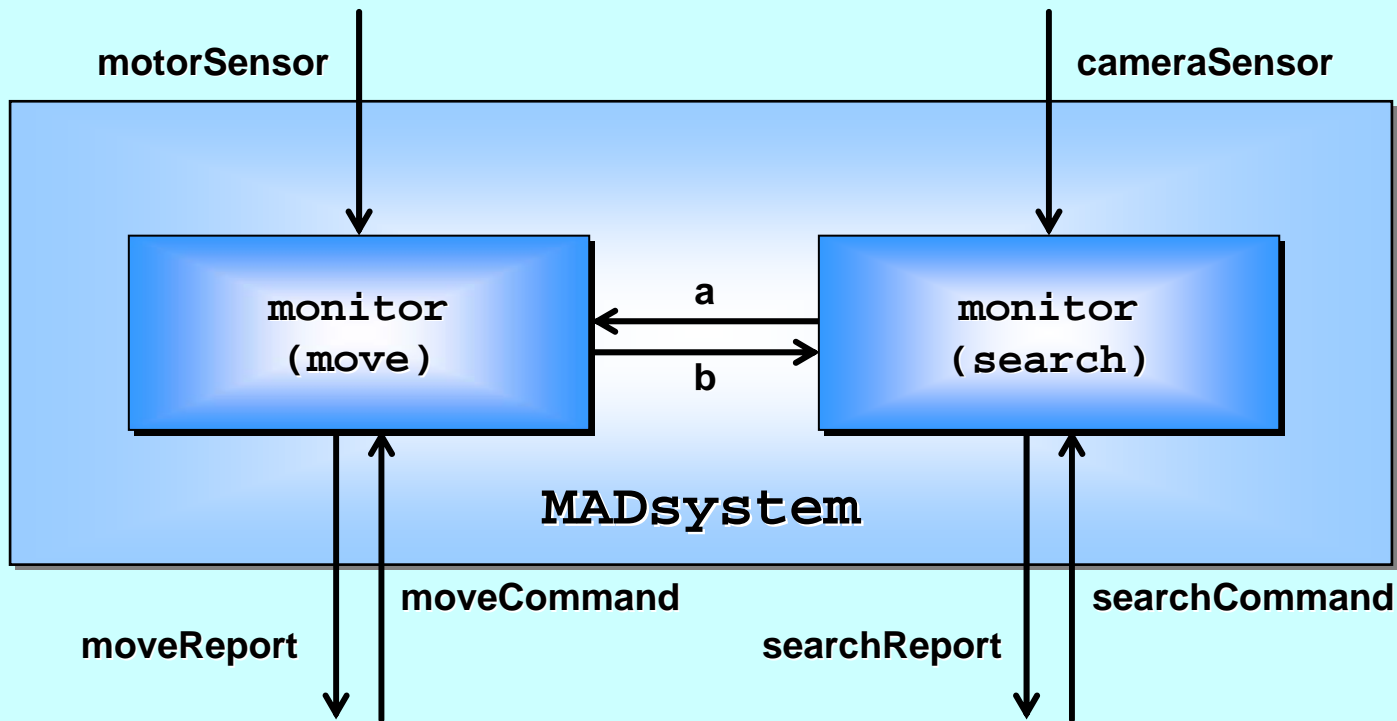
```
INT x:
sensor ? x
  SEQ
    ...  update local state with x (depends on mode)
    IF
      ...  task complete
        SEQ
          killYou ! kill
          report ! me
          running := FALSE
      TRUE
        SKIP
:
```

sensor

monitor
(mode)

killMe

killYou

report     command

<moveCommand, motorSensor,
searchCommand, cameraSensor>

Copyright  P.H.Welch, (2013)

**sensor**

```
monitor
(mode)
```

**killMe**

**killYou**

**report**   **command**

If the kill windows of the two monitors overlap, both will try to kill the other – resulting in deadlock.

**Kill Window**

```
INT x:
sensor ? x
  SEQ
    ...  update local state with x (depends on mode)
    IF
      ...  task complete
        SEQ
          killYou ! kill
          report ! me
          running := FALSE
      TRUE
        SKIP
  :
```

Copyright  P.H.Welch, (2013)

**average sensor data interval = 10 ms (randomised)**
**average sensor inputs per service = 100 (randomised)**
**➔**

**average service time = 1 second**
**kill window = 100 nanoseconds (approx.)**
**➔**

**chance of kill window overlap (deadlock) = 1/10,000,000**
**➔**

**time before 50% chance of deadlock = 90 days (approx.)**

**Kill Window**

```
INT x:
sensor ? x
  SEQ
    ...  update local state with x (depends on mode)
    IF
      ...  task complete
        SEQ
          killYou ! kill
          report ! me
          running := FALSE
      TRUE
        SKIP
  :
```

Copyright  P.H.Welch, (2013)

**chance of kill window overlap (deadlock) = 1/10,000,000**

> **This assumes each monitor runs on its own dedicated core …**
>
> **CCSP multicore scheduler dynamically bacthes processes to cores …**
>
> **If monitors are in the same batch, they will not deadlock …**

**chance of kill window overlap (deadlock) = 1/100,000,000 (approx.)**

➔

**time before 50% chance of deadlock = 2 years (approx.)**

**Kill Window**

```
INT x:
sensor ? x
  SEQ
    ...  update local state with x (depends on mode)
    IF
      ...  task complete
        SEQ
          killYou ! kill
          report ! me
          running := FALSE
      TRUE
        SKIP
  :
```

Copyright  P.H.Welch, (2013)

# Mutually Assured Destruction (revised implementation)

Communication between the monitors is mostly a *one-way* kill (from killer to killed).  Deadlock happens when both turn killer – *two-way* communications.

**Idea:** make communication between the monitors *always two-way* – either a **kill** in both directions (should both tasks complete around the same time) or a **kill** in one direction followed by an **ack** in the other (which will be most of the time).

**Claim:** this eliminates all deadlock (at the cost of an extra **ack**).

**INT**

sensor

**monitor
(mode)**

killMe

killYou

**PROTOCOL KILL
  CASE
    kill
:**

**INT**

command

report

**PROTOCOL REPORT
  CASE
    me          -- *task completed*
    she         -- *task abandoned*
:**

INT

sensor

**monitor'**
**(mode)**

killMe

killYou

INT

command

report

```
PROTOCOL KILL'
  CASE
    kill
    ack
:
```

```
PROTOCOL REPORT'
  CASE
    me        -- task completed
    she       -- task abandoned
    both      -- both completed
:
```

```
PROC monitor' (VAL INT mode, CHAN INT command?, sensor?,
               CHAN REPORT' report!,
               CHAN KILL' killYou!, killMe?)
  WHILE TRUE
    PRI ALT
      INT target:
      command ? target                 -- service requested
        service' (mode, target, sensor?, report!,
                  killYou!, killMe?)
      INT x:
      sensor ? x                        -- accept and discard
        SKIP
:
```

Copyright  P.H.Welch, (2013)

```
sensor │
       ▼
┌─────────────────────┐
│     monitor         │  ◄── killMe
│     (mode)          │  ──► killYou
└─────────────────────┘
    report │  ▲ command
           ▼  │
```

```occam
PROC service (VAL INT mode, target, CHAN INT sensor?,
              CHAN REPORT report!,
              CHAN KILL killYou!, killMe?)
  ...  local state and initialisation
  INITIAL BOOL running IS TRUE:
  WHILE running
    PRI ALT
      killMe ? kill
        ...  report 'she' and exit loop
      INT x:
      sensor ? x
        ...  process x
:
```

Copyright  P.H.Welch, (2013)

```
PROC service' (VAL INT mode, target, CHAN INT sensor?,
               CHAN REPORT' report!,
               CHAN KILL' killYou!, killMe?)
  ...  local state and initialisation
  INITIAL BOOL running IS TRUE:
  WHILE running
    PRI ALT
      killMe ? kill
        ... 'ack' the kill, report 'she' and exit loop
      INT x:
      sensor ? x
        ...  process x
:
```

**sensor**

```
monitor
(mode)
```

**killMe**

**killYou**

**report**   **command**

```
killMe ? kill
  SEQ
    report ! she
    running := FALSE
```

```
killMe ? kill
  SEQ
    killYou ! ack
    report ! she
    running := FALSE
```

Copyright  P.H.Welch, (2013)

```
PROC service' (VAL INT mode, target, CHAN INT sensor?,
               CHAN REPORT' report!,
               CHAN KILL' killYou!, killMe?)
  ...  local state and initialisation
  INITIAL BOOL running IS TRUE:
  WHILE running
    PRI ALT
      killMe ? kill
        ... 'ack' the kill, report 'she' and exit loop
      INT x:
      sensor ? x
        ...  process x
  :
```

Copyright  P.H.Welch, (2013)

sensor

monitor
(mode)

killMe

killYou

report    command

```
INT x:
sensor ? x
  SEQ
    ...  update local state with x (depends on mode)
    IF
      ...  task complete
        SEQ
          killYou ! kill
          report ! me
          running := FALSE
      TRUE
        SKIP
:
```

Copyright P.H.Welch, (2013)

sensor ↓

**monitor'**
**(mode)**

killMe ←
killYou →

report ↓ ↑ command

Each process knows what happened in the other – potentially a very useful side benefit.

```
INT x:
sensor ? x
  SEQ
    ...  update local state with x (depends on mode)
    IF
      ...  task complete
        SEQ
          PAR
            killYou ! kill        -- send and
            killMe ? CASE         -- receive in parallel
              ack
                report ! me
              kill
                report ! both
          running := FALSE
      TRUE
        SKIP
  :
```

Copyright  P.H.Welch, (2013)

sensor ↓

monitor'
(mode)

killMe ←
killYou →

report ↓ ↑ command

Key state information could easily be piggy-backed on the **kill** and **ack** signals – i.e. each process would know what the other found.

```
INT x:
sensor ? x
  SEQ
    ...   update local state with x (depends on mode)
    IF
      ...   task complete
        SEQ
          PAR
            killYou ! kill        -- send and
            killMe ? CASE         -- receive in parallel
              ack
                report ! me
              kill
                report ! both
          running := FALSE
      TRUE
        SKIP
  :
```

# Better ask for a model check …



**motorSensor**

**cameraSensor**

```
monitor'
(move)
```

a

b

```
monitor'
(search)
```

**MADsystem'**

**moveCommand**

**searchCommand**

**moveReport**

**searchReport**

**VERIFY DEADLOCK.FREE MADsystem'**  ✔  ☺ ☺ ☺

# Soak Testing



**motorSim**

**cameraSim**

**motorSensor**

**cameraSensor**

```
MADsystem'
```

**moveCommand**

**searchCommand**

**moveReport**

**searchReport**

**controllerSim**

**Only for confidence boosting – it will not deadlock
(assuming compiler, run-time kernel, microprocessor are OK)**

motorSensor

cameraSensor

**MADsystem'**

**moveCommand**

moveReport

searchReport

**searchCommand**

**This will not deadlock
(assuming compiler, run-time kernel, microprocessor are OK)**

# Mutually Assured Destruction
## (asynchronous channels?)

**motorSensor**

**cameraSensor**

**monitor (move)**

**a**

**b**

**monitor (search)**

**MADsystem**

**moveCommand**

**searchCommand**

**moveReport**

**searchReport**

If the channels were finitely buffered (with capacity greater than zero), the deadlock found with synchronous (i.e. zero-buffered) channels would not happen – both monitors would complete their kills, reports and service routines.

# Mutually Assured Destruction
## (asynchronous channels?)



If the channels were finitely buffered, deadlock is still possible – but less likely (exponentially) with increasing buffer size. Infinitely expandable buffer capacity would be needed to eliminate deadlock from the basic algorithm. For practical purposes, I would feel safe with a capacity of 3.

# Mutually Assured Destruction (asynchronous channels?)



However, there is a nasty problem. If both monitors send a **kill**, neither is taken and they remain lurking in the buffered channels. Some time in the next service cycle, both will strike and the services will be erroneously aborted.

# Mutually Assured Destruction (asynchronous channels?)

**motorSensor**

**cameraSensor**

| monitor (move) | ← **a** → | monitor (search) |

**b**

**MADsystem**

**moveCommand**

**searchCommand**

**moveReport**

**searchReport**

This could be overcome by counting cycles and *sequence numbering* the `kill` signals: just ignore any `kill` with a number less than the current count. This adds complexity and run-time overhead.

# Mutually Assured Destruction (asynchronous channels?)

**motorSensor**

**cameraSensor**

```
         monitor          a          monitor
         (move)      <----------     (search)
                     ---------->
                          b
```

**MADsystem**

**moveCommand**

**searchCommand**

**moveReport**

**searchReport**

This could be overcome by counting cycles and *sequence numbering* the `kill` signals: just ignore any `kill` with a number less than the current count. Further, this only works if the processes engaged in MAD are in lock-step (which they are in this scenario, but not in general).

# Mutually Assured Destruction
## (asynchronous channels?)



Alternatively, the mess could be sorted out by the **Controller** process. When/if it gets two **me** reports from the monitors, it tells each monitor (as part of its next command) to read and discard an incoming **kill**. Again, this adds complexity – we shouldn't have a mess to clean up!

# Mutually Assured Destruction (asynchronous channels?)

**motorSensor**

**cameraSensor**

```
monitor
(move)
```

a

b

```
monitor
(search)
```

**MADsystem**

**moveCommand**

**searchCommand**

**moveReport**

**searchReport**

Alternatively, the mess could be sorted out by the **Controller** process. When/if it gets two **me** reports from the monitors, it tells each monitor (as part of its next command) to read and discard an incoming **kill**. Further, this assumes a **Controller**, which processes engaged in **MAD** may not have.

# The Joy of Sync

Process oriented design ...

Synchronous communications …

Synchronous barriers …

Mutually assured destruction …

Non-blocking barriers …

Performance …

# Non-Blocking Barriers

Recently (2012) introduced to MPI, *non-blocking barrier synchronisation* seems, at first glance, a contradiction of terms … the whole point of a *barrier* is to *block* until all parties are there!

When we have completed our work before a *barrier*, we normally synchronise on it – thereby notifying that we are there and waiting for the others.

Recall …

# Barriers

Processes may synchronise on more than one barrier:



```
BARRIER b, c:
PAR i = 0 FOR n ENROLL b, c
  worker (i, b, c)
```

To synchronise on a barrier:

```
SYNC b
```
or
```
SYNC c
```

# Barriers

Barriers are commonly used to synchronise multiple *phases* of computation between a set of processes.  Within each phase, other synchronisations (channel / barrier) may take place:

```
PROC worker (VAL INT id, BARRIER b, c)
  ...   local declarations / initialisation
  WHILE running
    SEQ
      SYNC b
      ...   observe neighbourhood phase
      SYNC c
      ...   change neighbourhood phase
:
```

**All workers do this together – all see the same thing …**

**All workers do this together – may need to negotiate …**

# Non-Blocking Barriers

Recently (2012) introduced to MPI, *non-blocking barrier synchronisation* seems, at first glance, a contradiction of terms … the whole point of a *barrier* is to *block* until all parties are there!

When we have completed our work before a *barrier*, we normally synchronise on it – thereby notifying that we are there and waiting for the others.

However, if there is something we can be getting on with that does not disturb our fellow *synchronisers,* (e.g. preparatory work for the phase following the *barrier*), it would be good to be able to do so. Only when we need stuff that depends on the other *synchronisers,* should we have to wait for them.

# Blocking Barrier Sync (MPI)

```
...   phase 0 computation
MPI_Barrier (b, ...);              // wait for everyone ...
...   preparatory work for next phase
...   phase 1 computation
```

# Blocking Barrier Sync (occam-π)

```
SEQ
  ...   phase 0 computation
  SYNC b                           -- wait for everyone ...
  ...   preparatory work for next phase
  ...   phase 1 computation
```

Copyright P.H.Welch, (2013)

# Non-Blocking Barrier Sync (MPI)

```
...   phase 0 computation
MPI_IBarrier (b, ...);          // hey, I'm done ...
...   preparatory work for next phase
MPI_WBarrier (b, ...);          // I'm waiting now ...
...   phase 1 computation
```

# Non-Blocking Barrier Sync (o-π)

```
SEQ
  ...   phase 0 computation
  PAR
    SYNC b                      -- hey, I'm done ...
    ...   preparatory work for next phase
  ...  phase 1 computation
```

# Non-Blocking Barrier Sync (0-π)

```
SEQ
  ...   phase 0 computation
  PAR
    SYNC b                        -- hey, I'm done ...
    ...   preparatory work for next phase
  ...   phase 1 computation
```

The **SYNC** registers that we have arrived at the barrier and lets all move forward when the rest arrive. In parallel with the above, we get on with our *preparatory work*.  ☺

When our *preparatory work* is complete, if all the others have reached the barrier, the **SYNC** will have completed – so the **PAR** completes and we immediately move on to **phase 1**.  And we have not delayed the others.  ☺

When our *preparatory work* is complete, if not all the others have reached the barrier, the **SYNC** will not have completed. We wait for the others at the **SYNC** before moving on to **phase 1** – as we must!  ☺

Copyright  P.H.Welch, (2013)

# Non-Blocking Barrier Sync (o-π)

```
SEQ
  ...   phase 0 computation
  PAR
    SYNC b                        -- hey, I'm done ...
    ...  preparatory work for next phase
  ...   phase 1 computation
```

The **SYNC** registers that we have arrived at the barrier and lets all move forward when the rest arrive. In parallel with the above, we get on with our *preparatory work*.  ☺

When our *preparatory work* is co̶m̶p̶l̶e̶t̶e̶ ̶i̶f̶ ̶t̶h̶e̶ others have reached the barrier, the **SYNC** will h̶a̶v̶e̶ ̶c̶o̶m̶p̶l̶e̶t̶e̶d̶, so the **PAR** completes and we immediately mov̶e̶ ̶t̶o̶ ̶p̶h̶a̶s̶e̶ ̶1̶.  And we have not delayed the others.  ☺

Wh̶e̶n̶ ̶t̶h̶e̶ *preparatory work* is complete, if not all the others have reached the barrier, the **SYNC** will not have completed. We wait for the others at the **SYNC** before moving on to **phase 1** – as we must!  ☺

**Nothing new in occam-π is needed for this.**

# The Joy of Sync

Process oriented design ...

Synchronous communications …

Synchronous barriers …

Mutually assured destruction …

Non-blocking barriers …

Performance …

# Performance

Take a ring of $n$ processes …

Each process connects to all …

In parallel, each process sends and receives $m$ messages (e.g. its *id number*) to all, including itself …

That's $mn^2$ messages …

How long per message?

**P(n-1)**

**P(0)**

**P(1)**

**P(2)**

**P(3)**

# Performance

Andrew's "Say Hello to Everyone" Benchmark (v1) \*

$2n^2$ **processes**

$mn^2$ **messages**

```
[n][n]CHAN INT c:

PAR i = 0 FOR n
  PAR
    PAR j = 0 FOR n        -- for each j in parallel,
      SEQ k = 0 FOR m      -- send m messages (i to j)
        c[i][j] ! i
    PAR j = 0 FOR n        -- for each j in parallel,
      SEQ k = 0 FOR m      -- receive m messages (j to i)
        INT x:
        SEQ
          c[j][i] ? x
          ASSERT (x = j)   -- sanity check
```

# Performance

Andrew's "Say Hello to Everyone" Benchmark (v2) *

**2n processes**

**$mn^2$ messages**

```
[n][n]CHAN INT c:

PAR i = 0 FOR n
  PAR
    SEQ j = 0 FOR n       -- for each j in series,
      SEQ k = 0 FOR m       -- send m messages (i to j)
        c[i][j] ! i
    SEQ j = 0 FOR n       -- for each j in series,
      SEQ k = 0 FOR m       -- receive m messages (j to i)
        INT x:
        SEQ
          c[j][i] ? x
          ASSERT (x = j)    -- sanity check
```

# Performance

P(n-1)

P(0)

P(1)

Take a ring of **N** processes …

Each

In para
and re
its id)

The following observations were made using KRoC 1.5.0-pre5, Ubuntu 11.04 (natty) on an Intel i7 quad-core processor with hyperthreading (i.e. 8 virtual cores), running at 2 MHz.

The benchmark timings were averaged from 10 separate runs, with negligible variance.

That's **mn²** messages …

How long per message?

P(3)

# Performance

**P(n-1)**

**P(0)**

Take a ring of **N** processes...

**1 core / 20,000 nodes / v2**

| Monitor | Edit | View | Help |

System | Processes | Resources | File Systems

**CPU History**

- CPU1 11.9%
- CPU2 0.0%
- CPU3 44.4%
- CPU4 57.4%
- CPU5 4.0%
- CPU6 5.0%
- CPU7 4.8%
- CPU8 1.0%

**Memory and Swap History**

Memory
6.2 GiB (80.9%) of 7.7 GiB

Swap
0 bytes (0.0%) of 7.9 GiB

**Network History**

Receiving    179 bytes/s         Sending      0 bytes/s
Total Received   42.9 MiB        Total Sent   10.8 MiB

# Performance

**P(n-1)**

**P(0)**

Take a ring of N processes ...

E ...

In ...
an ...
its ...

T ...

H ...

**2 cores / 20,000 nodes / v2**

Monitor   Edit   View   Help

System | Processes | **Resources** | File Systems

**CPU History**

100 %
50 %
0 %

60 seconds        50        40        30        20        10        0

CPU1 95.0%          CPU2 5.0%          CPU3 2.0%          CPU4 3.0%
CPU5 100.0%         CPU6 0.0%          CPU7 7.0%          CPU8 15.0%

**Memory and Swap History**

100 %
50 %
0 %

60 seconds        50        40        30        20        10        0

Memory                              Swap
6.2 GiB (80.8%) of 7.7 GiB          0 bytes (0.0%) of 7.9 GiB

**Network History**

8.0 KiB/s
4.0 KiB/s
0.0 KiB/s

60 seconds        50        40        30        20        10        0

Receiving                Sending
Total Received           Total Sent
42.8 MiB    0 bytes/s    9.6 MiB    227 bytes/s

**P(2)**

# Performance

**P(n-1)**

**P(0)**

Take a ring of N processes

E...

In...
a...
its...

T...

H...



**4 cores / 20,000 nodes / v2**

Monitor  Edit  View  Help

System | Processes | Resources | File Systems

**CPU History**

100 %
50 %
0 %
60 seconds    50    40    30    20    10    0

CPU1 17.8%    CPU2 11.9%    CPU3 100.0%    CPU4 100.0%
CPU5 14.9%    CPU6 81.9%    CPU7 0.9%    CPU8 100.0%

**Memory and Swap History**

100 %
50 %
0 %
60 seconds    50    40    30    20    10    0

Memory                          Swap
6.2 GiB (80.7%) of 7.7 GiB      0 bytes (0.0%) of 7.9 GiB

**Network History**

8.0 KiB/s
4.0 KiB/s
0.0 KiB/s
60 seconds    50    40    30    20    10    0

Receiving          356 bytes/s        Sending            0 bytes/s
Total Received     42.7 MiB           Total Sent         9.6 MiB

# Performance

**P(n-1)**

**P(0)**

Take a ring of N processes.

E...

In...
a...
its...

T...

H...

**8 cores / 20,000 nodes / v2**

Monitor   Edit   View   Help

System | Processes | Resources | File Systems

**CPU History**

100 %
50 %
0 %
60 seconds        50        40        30        20        10        0

CPU1 100.0%   CPU2 100.0%   CPU3 100.0%   CPU4 100.0%
CPU5 100.0%   CPU6 100.0%   CPU7 99.0%    CPU8 100.0%

**Memory and Swap History**

100 %
50 %
0 %
60 seconds        50        40        30        20        10        0

Memory
6.2 GiB (80.7%) of 7.7 GiB

Swap
0 bytes (0.0%) of 7.9 GiB

**Network History**

8.0 KiB/s
4.0 KiB/s
0.0 KiB/s
60 seconds        50        40        30        20        10        0

Receiving          0 bytes/s          Sending          0 bytes/s
Total Received     42.6 MiB           Total Sent        9.5 MiB

# Performance

P(n-1)

P(0)

P(1)

P(2)

P(3)

So … how long per message?

6,000 nodes …

(V2) 12,000 processes …

30 messages node-to-node …

1.08 billion messages …

8 cores …

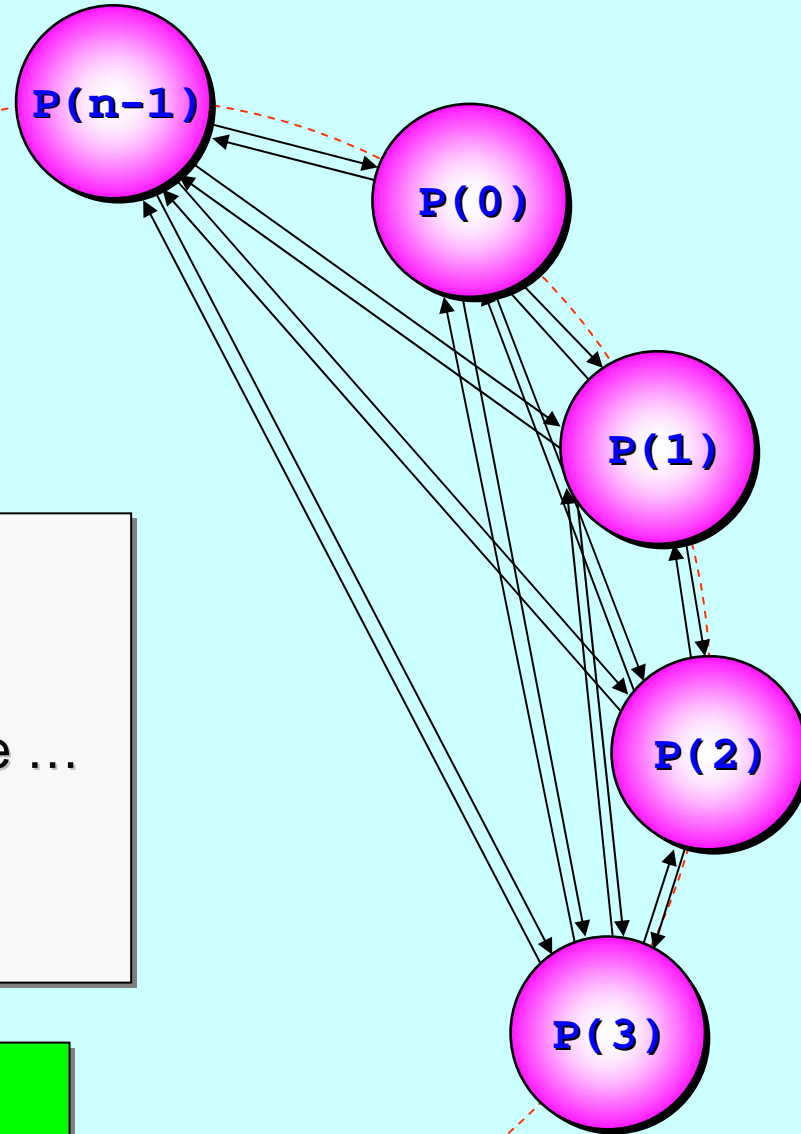7 nanoseconds

# Performance

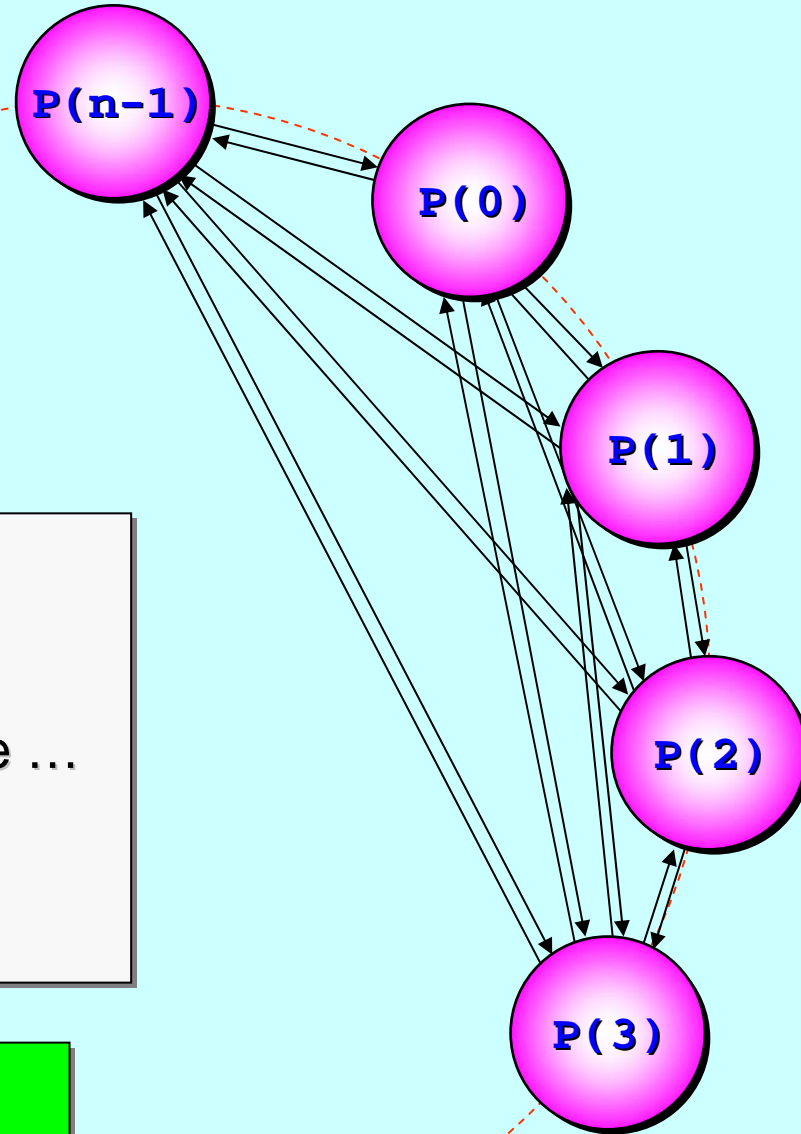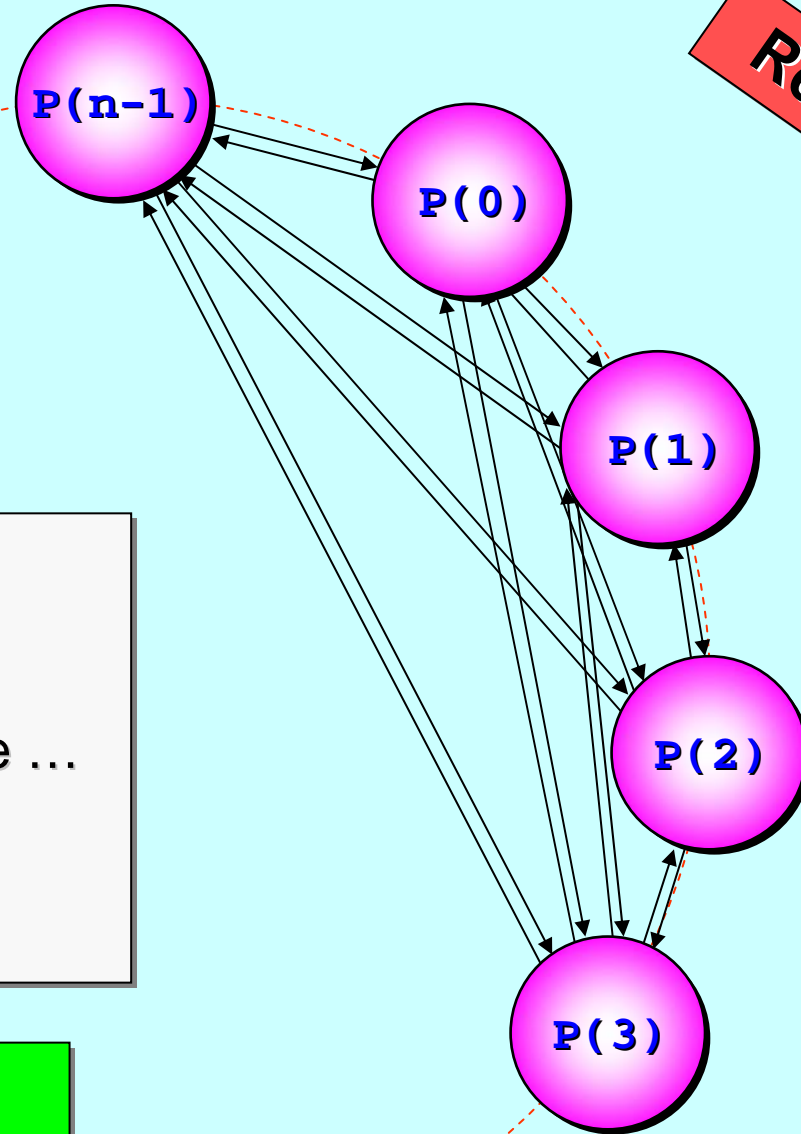So … how long per message?

10,000 nodes …

(V2) 20,000 processes …

30 messages node-to-node …

3 billion messages …

8 cores …

## 7 nanoseconds

P(n-1)

P(0)

P(1)

P(2)

P(3)

# Performance

**P(n-1)**

**P(0)**

**P(1)**

**P(2)**

**P(3)**

So … how long per message?

20,000 nodes …

(V2) 40,000 processes …

30 messages node-to-node …

12 billion messages …

8 cores …

8 **nanoseconds**
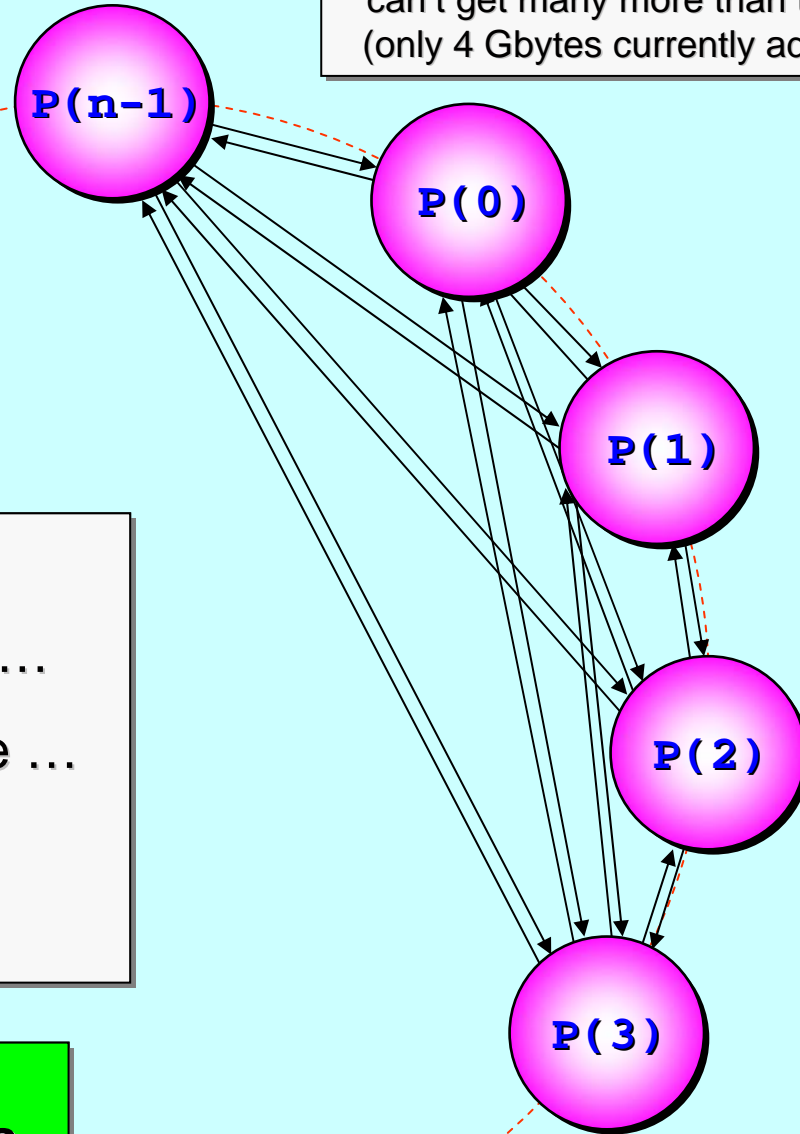
# Performance

P(n-1)

P(0)

P(1)

P(2)

P(3)

**Repeat.**

So … how long per message?

6,000 nodes …

(V2) 12,000 processes …

30 messages node-to-node …

1.08 billion messages …

8 cores …

7 **nanoseconds**

# **Performance**

**P(n-1)**

**P(0)**

**P(1)**

**P(2)**

**P(3)**

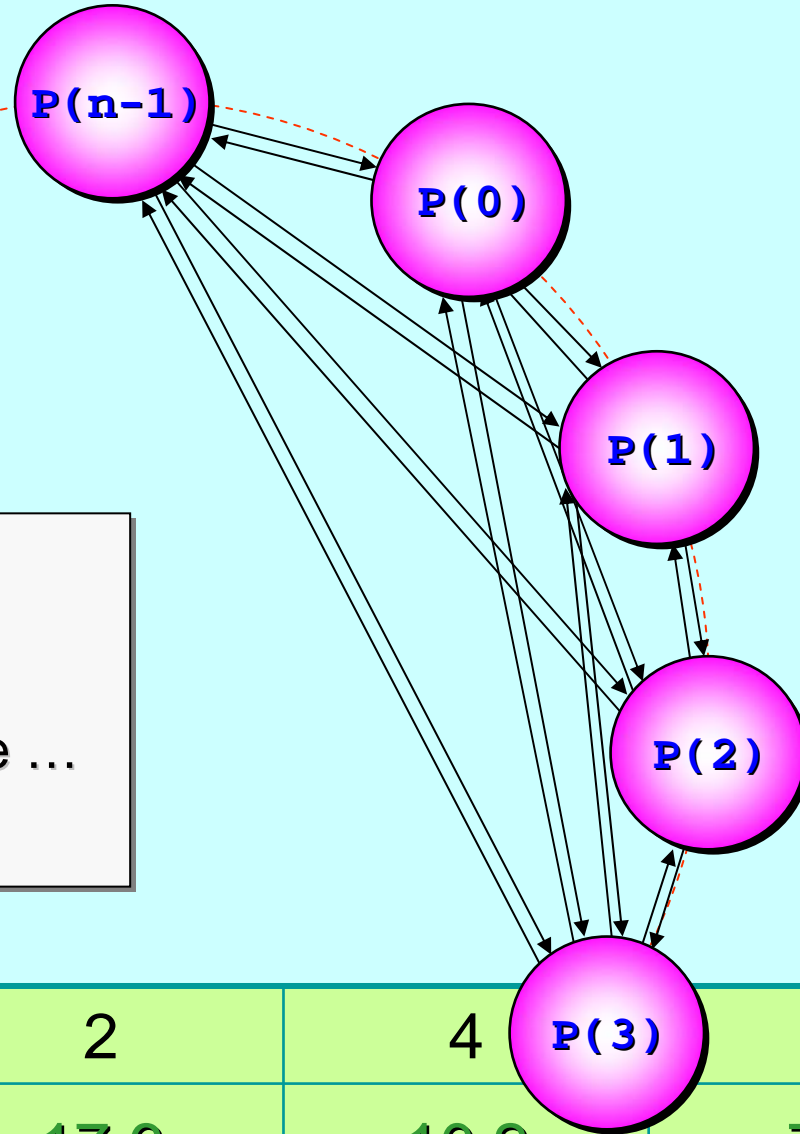* can't get many more than this
  (only 4 Gbytes currently addressable).

So … how long per message?

6,000 nodes …

(V1) 72 million processes* …

30 messages node-to-node …

1.08 billion messages …

8 cores …

## 14 nanoseconds

# Performance

**P(n-1)**

**P(0)**

**P(1)**

**P(2)**

**P(3)**

So … how long per message?

6,000 nodes …

(V2) 12,000 processes …

30 messages node-to-node …

1.08 billion messages …

| # cores | 1 | 2 | 4 | 8 |
|---|---|---|---|---|
| per message (nanoseconds) | 34.5 | 17.6 | 10.2 | 7.0 |
| speed up | 1.0 | 2.0 | 3.4 | 4.9 |

# **Performance**
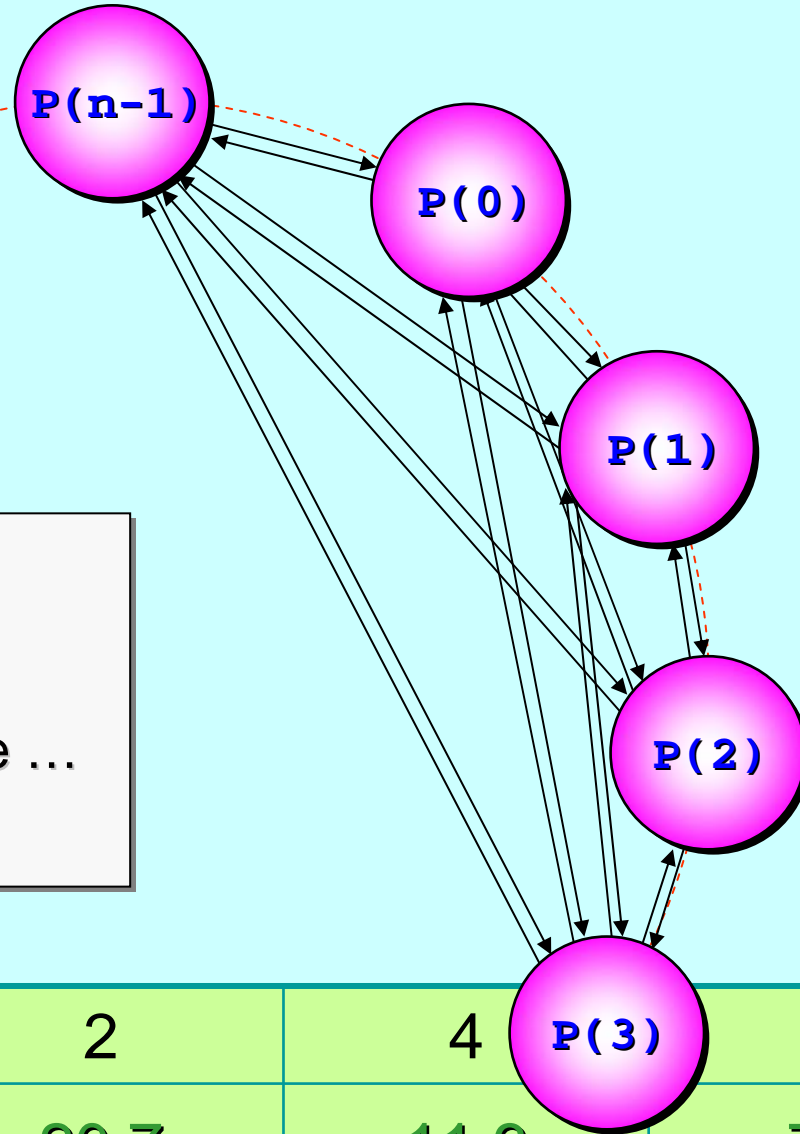
**P(n-1)**

**P(0)**

**P(1)**

**P(2)**

**P(3)**

So … how long per message?

10,000 nodes …

(V2) 20,000 processes …

30 messages node-to-node …

3 billion messages …

| # cores | 1 | 2 | 4 | 8 |
|---|---|---|---|---|
| per message (nanoseconds) | 37.5 | 17.4 | 10.3 | 7.2 |
| speed up | 1.0 | 2.2 | 3.6 | 5.2 |

# Performance

**P(n-1)**

**P(0)**

**P(1)**

**P(2)**

**P(3)**

So … how long per message?

20,000 nodes …

(V2) 40,000 processes …

30 messages node-to-node …

12 billion messages …

| # cores | 1 | 2 | 4 | 8 |
|---|---|---|---|---|
| per message (nanoseconds) | 61.8 | 20.7 | 11.3 | 7.7 |
| speed up | 1.0 | 3.0 | 5.5 | 8.0 |

# **Performance**

**P(n-1)**

**P(0)**

**P(1)**

**P(2)**

**P(3)**

**Repeat.**

So … how long per message?

6,000 nodes …

(V2) 12,000 processes …
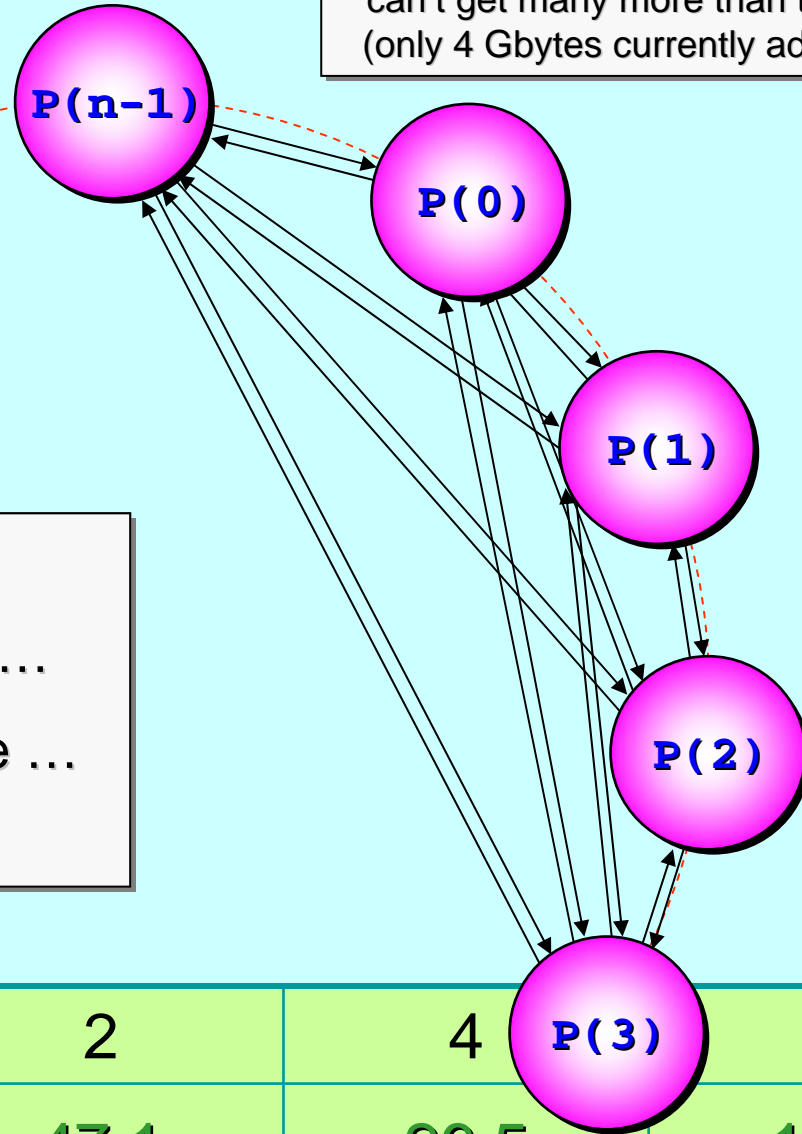
30 messages node-to-node …

1.08 billion messages …

| # cores | 1 | 2 | 4 | 8 |
|---|---|---|---|---|
| per message (nanoseconds) | 34.5 | 17.6 | 10.2 | 7.0 |
| speed up | 1.0 | 2.0 | 3.4 | 4.9 |

# **Performance**

**P(n-1)**

**P(0)**

**P(1)**

**P(2)**

**P(3)**

So … how long per message?

6,000 nodes …

(V1) 72 million processes* …

30 messages node-to-node …

1.08 billion messages …

| # cores | 1 | 2 | 4 | 8 |
|---|---|---|---|---|
| per message (nanoseconds) | 120.0 | 47.1 | 26.5 | 14.0 |
| speed up | 1.0 | 2.5 | 4.5 | 8.6 |

# The Joy of Sync

Process oriented design ...

Synchronous communications …

Synchronous barriers …

Mutually assured destruction …

Non-blocking barriers …

Performance …

Any questions?