

# Virtual Synchrony, Paxos, and Beyond

Ken Birman

*Cornell University. CS5410 Fall 2008.*





# Virtual Synchrony

- A powerful programming model!
- Key idea: equate “group” with “data abstraction”
  - Each group implements some object
  - An application can belong to many groups
- *Virtual synchrony* is the associated consistency model:
  - Process groups with state transfer, automated fault detection and membership reporting
  - Ordered reliable multicast, in several flavors
  - Extremely good performance
- Assumes that the membership oracle is available



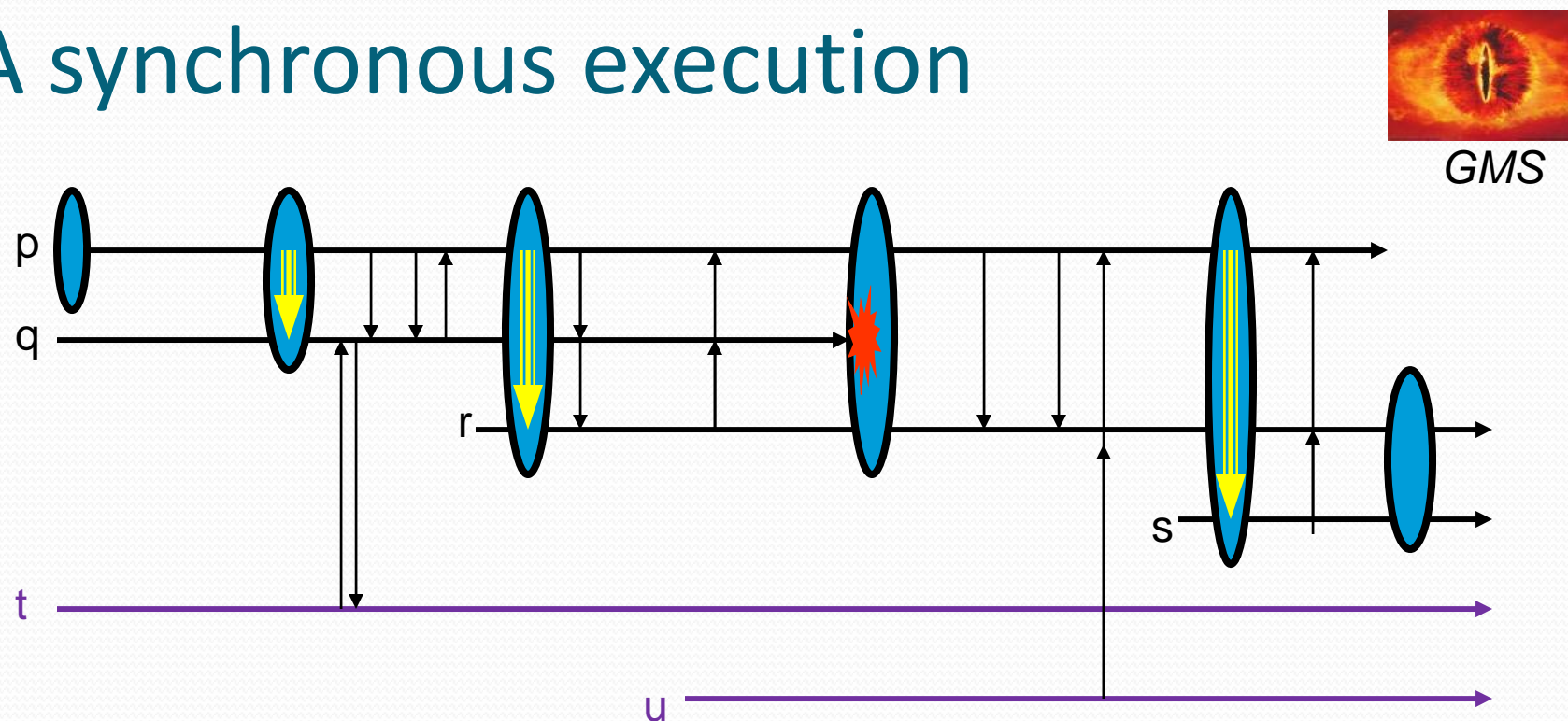


# Why “virtual” synchrony?

- What would a synchronous execution look like?
- In what ways is a “virtual” synchrony execution not the same thing?



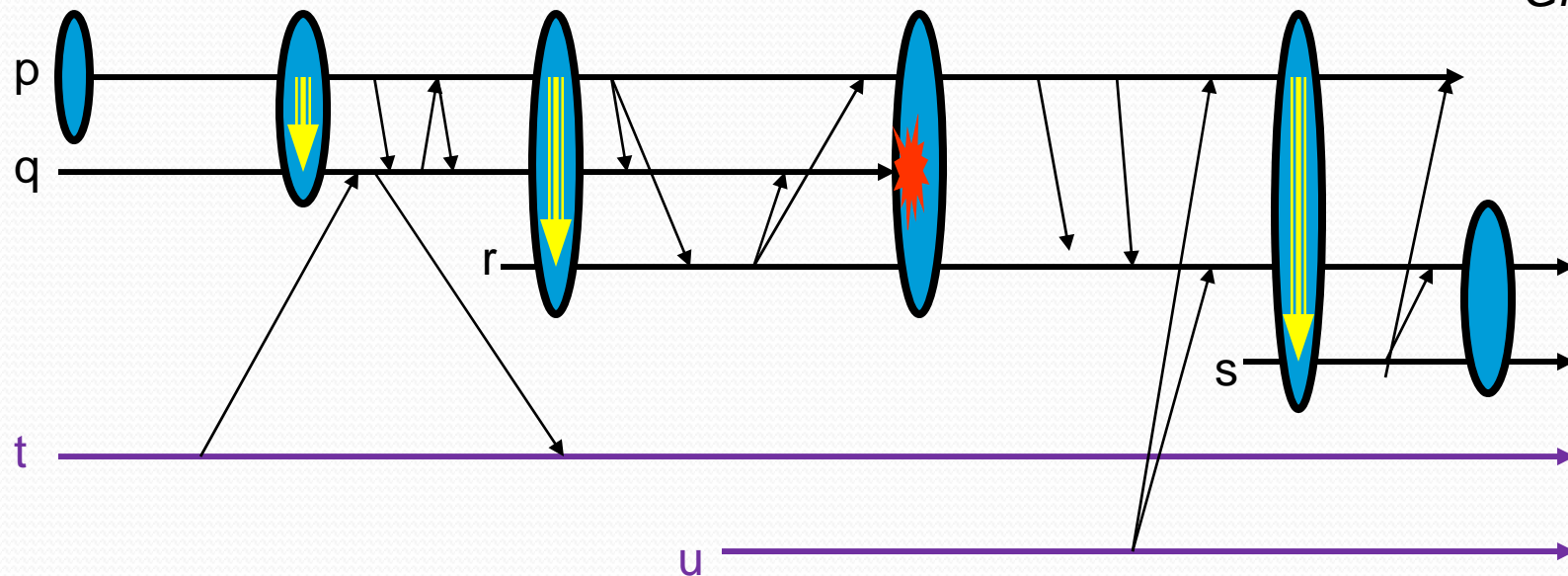
# A synchronous execution



- With *true* synchrony executions run in genuine lock-step.
- GMS treated as a membership oracle



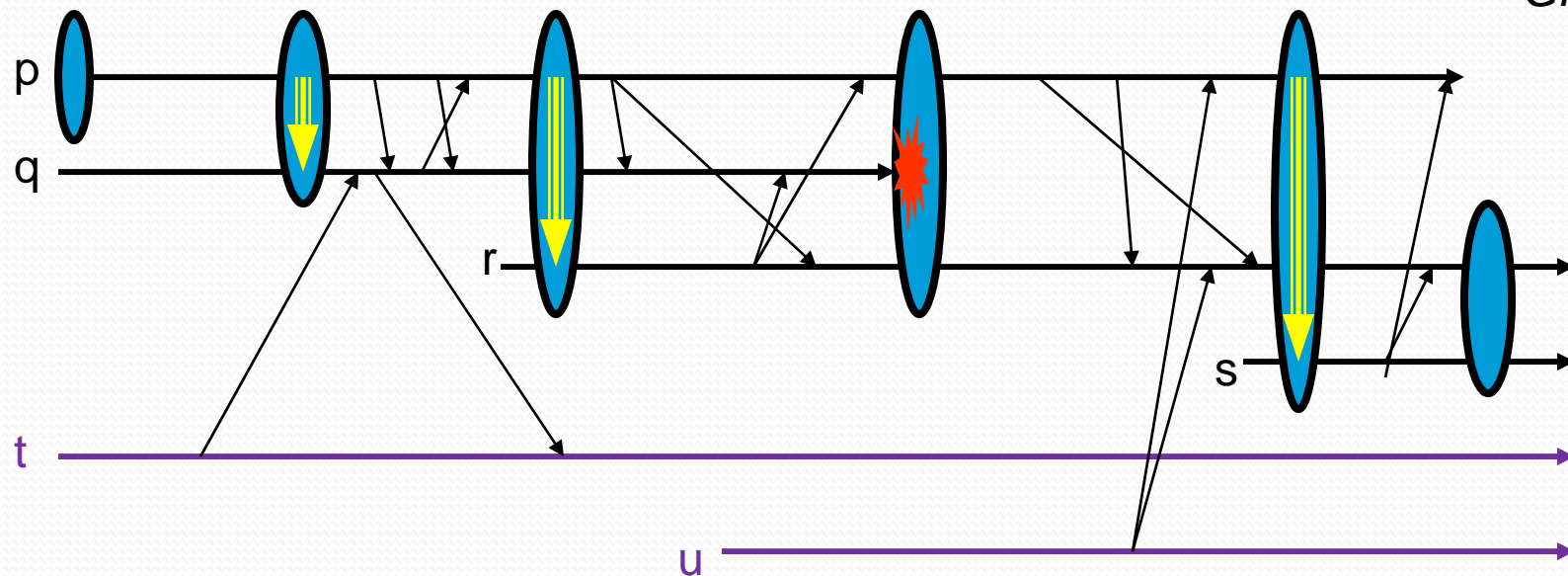
# Virtual Synchrony at a glance



- With *virtual* synchrony executions only look “lock step” to the application



# Virtual Synchrony at a glance



- We use the weakest (hence fastest) form of communication possible





## f/c/abcast, flush...

- Last week we saw how the GMS could support protocols that have various forms of ordering
  - One option is to run them “inside” the GMS
- But suppose that we use the GMS to manage membership of processes other than the GMS servers
  - *This is easy do to! Our protocols didn’t “need” to run inside the GMS per-se!*
  - Group runs “outboard” and protocols send messages directly between the members
  - GMS involved only if the group view changes



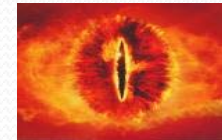


# Chances to “weaken” ordering

- For example, imagine a data replication object with variables “inside”, and suppose that any conflicting updates are synchronized using some form of locking
  - Multicast sender will have mutual exclusion
  - Hence simply because we used locks, cbcast delivers conflicting updates in order they were performed!
- If our system ever *does* see concurrent multicasts... they must not have conflicted. So it won't matter if cbcast delivers them in different orders at different recipients!

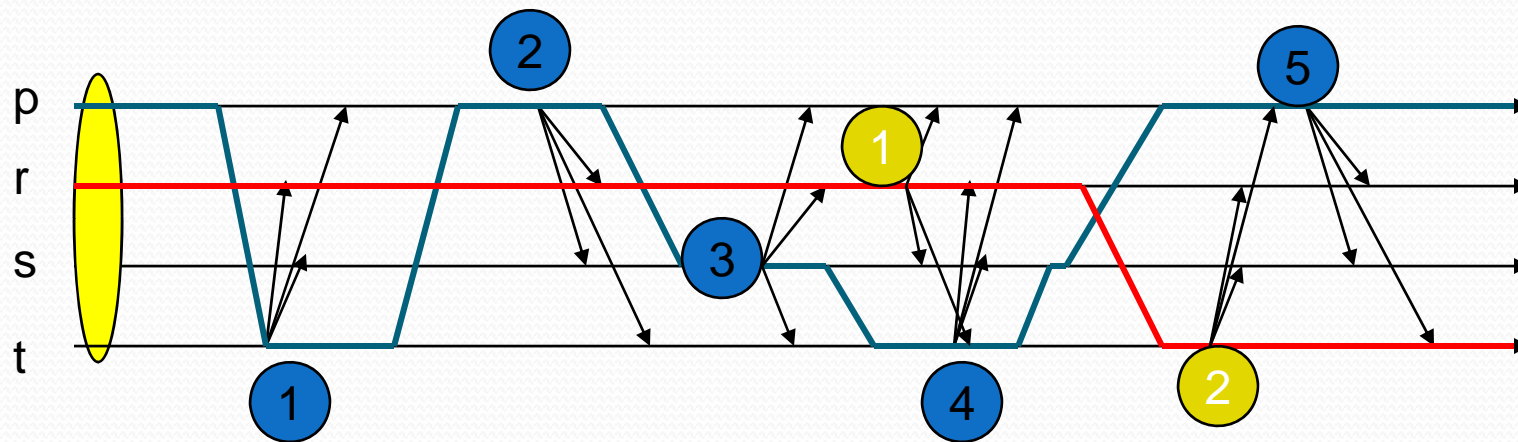


# Causally ordered updates



GMS

- Each thread corresponds to a different lock



- The “group” is best visualized as a kind of object that has a replica associated with each process
- Within it, red “events” never conflict with blue!





# Strong to weak...

- Definition: a multicast is *safe*, also called *dynamically uniform*, if:
  - If any group member delivers it to the application layer, then every group member will do so, unless it fails
  - And this is true even if the first deliveries are in processes that fail
- Our *fbcast* and *cbcast* and *abcast* protocols from last week were not safe!
  - They delivered “early” and hence if the first processes to deliver copies failed, the event might be lost





# Strong to weak

- A safe/dynamically uniform protocol needs two phases
  - Phase 1: Deliver a copy to at least a majority of the members of the current view. Recipients acknowledge
  - ... but instead of delivering to the application layer, they retain these copies in buffers
  - Phase 2: Sender detects that a majority have a copy, tells recipients that now they can deliver
- Much slower... but now a failure can't erase a multicast.
- Same ordering options exist





# In general?

- Start by thinking in terms of safe abcast (== gbcast!)
- Then replace “safe” (dynamic uniformity) protocols with a standard multicast when possible
- Weaken ordering, replacing abcast with cbcast
- Weaken even more, replacing cbcast with fbcast





# More tricks of the trade

- Multicast primitives usually can support replies!
  - No replies: the multicast is *asynchronous*
  - One reply: like an *anycast* – all receive, but any one reply will suffice (first one wins if several reply...)
  - $C$  replies: for a constant  $C$  (rarely supported)
  - *ALL* replies: wait for every group member to reply
- *Failure*: treated as a “null reply”
- Want speed? Ask for as few replies as possible!





# Why “virtual” synchrony?

- The user sees what looks like a synchronous execution
  - Simplifies the developer’s task
- But the actual execution is rather concurrent and asynchronous
  - Maximizes performance
  - Reduces risk that lock-step execution will trigger correlated failures





# Correlated failures

- Observation: virtual synchrony makes these less likely!
  - Recall that many programs are buggy
  - Often these are Heisenbugs (order sensitive)
- With lock-step execution each group member sees group events in identical order
  - So all die in unison
- With virtual synchrony orders differ
  - So an order-sensitive bug might only kill one group member!





# Programming with groups

- Many systems just have one group
  - E.g. replicated bank servers
  - Cluster mimics one highly reliable server
- But we can also use groups at finer granularity
  - E.g. to replicate a shared data structure
  - Now one process might belong to many groups
- A further reason that different processes might see different inputs and event orders





# Embedding groups into “tools”

- We can design a groups API:
  - `pg_join()`, `pg_leave()`, `cbcast()`...
- But we can also use groups to build other higher level mechanisms
  - Distributed algorithms, like snapshot
  - Fault-tolerant request execution
  - Publish-subscribe





# Distributed algorithms

- Processes that might participate join an appropriate group
- Now the group view gives a simple leader election rule
  - Everyone sees the same members, in the same order, ranked by when they joined
  - Leader can be, e.g., the “oldest” process





# Distributed algorithms

- A group can easily solve consensus
  - Leader multicasts: “what’s your input”?
  - All reply: “Mine is 0. Mine is 1”
  - Initiator picks the most common value and multicasts that: the “decision value”
  - If the leader fails, the new leader just restarts the algorithm
- Puzzle: Does FLP apply here?





# Distributed algorithms

- A group can easily do consistent snapshot algorithm
  - Either use cbcast throughout system, or build the algorithm over gbcast
  - Two phases:
    - Start snapshot: a first cbcast
    - Finished: a second cbcast, collect process states and channel logs





# Distributed algorithms: Summary

- Leader election
- Consensus and other forms of agreement like voting
- Snapshots, hence deadlock detection, auditing, load balancing





# More tools: fault-tolerance

- Suppose that we want to offer clients “fault-tolerant request execution”
  - We can replace a traditional service with a group of members
  - Each request is assigned to a primary (ideally, spread the work around) and a backup
    - Primary sends a “cc” of the response to the request to the backup
  - Backup keeps a copy of the request and steps in only if the primary crashes before replying
- Sometimes called “coordinator/cohort” just to distinguish from “primary/backup”





# Publish / Subscribe

- Goal is to support a simple API:
  - Publish(“topic”, message)
  - Subscribe(“topic”, event\_handler)
- We can just create a group for each topic
  - Publish multicasts to the group
  - Subscribers are the members





# Scalability warnings!

- Many existing group communication systems don't scale incredibly well
  - E.g. JGroups, Ensemble, Spread
  - Group sizes limited to perhaps 50-75 members
  - And individual processes limited to joining perhaps 50-75 groups (Spread: see next slide)
- Overheads soar as these sizes increase
  - Each group runs protocols oblivious of the others, and this creates huge inefficiency





# Publish / Subscribe issue?

- We could have thousands of topics!
  - Too many to directly map topics to groups
- Instead map topics to a *smaller set* of groups.
  - SPREAD system calls these “lightweight” groups
  - Mapping will result in inaccuracies... Filter incoming messages to discard any not actually destined to the receiver process
- Cornell’s new QuickSilver system will instead directly support immense numbers of groups





# Other “toolkit” ideas

- We could embed group communication into a framework in a “transparent” way
  - Example: CORBA fault-tolerance specification does lock-step replication of deterministic components
  - The client simply can’t see failures
    - But the determinism assumption is painful, and users have been unenthusiastic
    - And exposed to correlated crashes





# Other similar ideas

- There was some work on embedding groups into programming languages
  - But many applications want to use them to link programs coded in different languages and systems
  - Hence an interesting curiosity but just a curiosity
- More work is needed on the whole issue



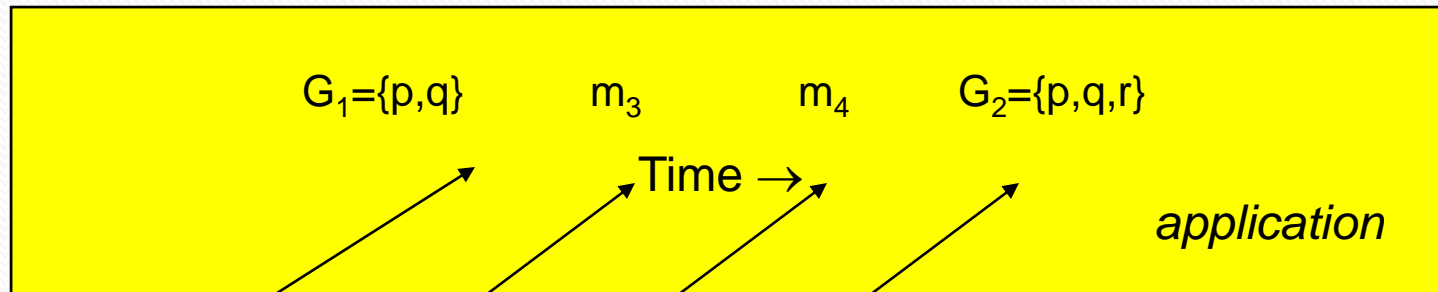


# Existing toolkits: challenges

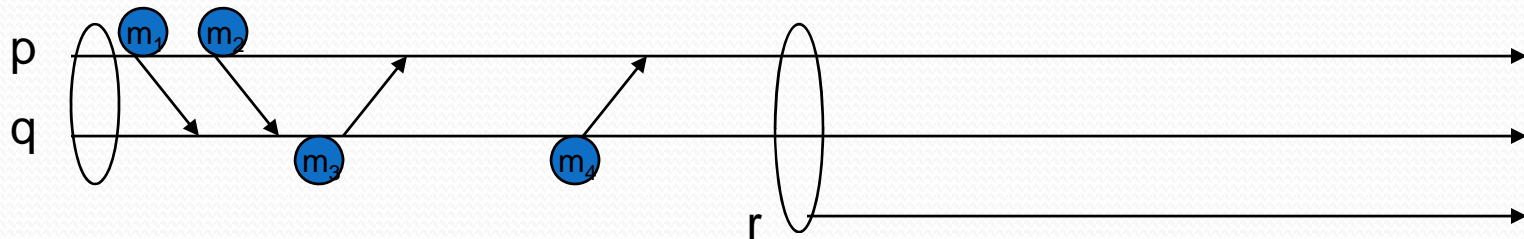
- Tensions between threading and ordering
  - We need concurrency (threads) for perf.
  - Yet we need to preserve the order in which “events” are delivered
- This poses a difficult balance for the developers



# Preserving order



*Group Communication Subsystem: A library linked to the application, perhaps with its own daemon processes*







# The tradeoff

- If we deliver these upcalls in separate threads, concurrency increases but order could be lost
- If we deliver them as a list of event, application receives events in order... but if it uses thread pools (most famous version of this is called SEDA), the order is lost





# Solution used in Horus

- This system
  - Delivered upcalls using an event model
  - Each event was numbered
  - User was free to
    - Run a single-threaded app
    - Use a SEDA model
- Toolkit included an “enter/leave region in order” synchronization primitive
  - Forced threads to enter in event-number order





# Other toolkit “issues”

- Does the toolkit distinguish members of a group from clients of that group?
  - In Isis system, a client of a group was able to multicast to it, with vsync properties
  - But only members received events
- Does the system offer properties “across group boundaries”?
  - For example, using cbcast in multiple groups



# Features of major virtual synchrony platforms

- Isis: First and no longer widely used
  - But was perhaps the most successful; has major roles in NYSE, Swiss Exchange, French Air Traffic Control system (two major subsystems of it), US AEGIS Naval warship
  - Also was first to offer a publish-subscribe interface that mapped topics to groups



# Features of major virtual synchrony platforms

- Totem and Transis
  - Sibling projects, shortly after Isis
  - Totem (UCSB) went on to become Eternal and was the basis of the CORBA fault-tolerance standard
  - Transis (Hebrew University) became a specialist in tolerating partitioning failures, then explored link between vsync and FLP

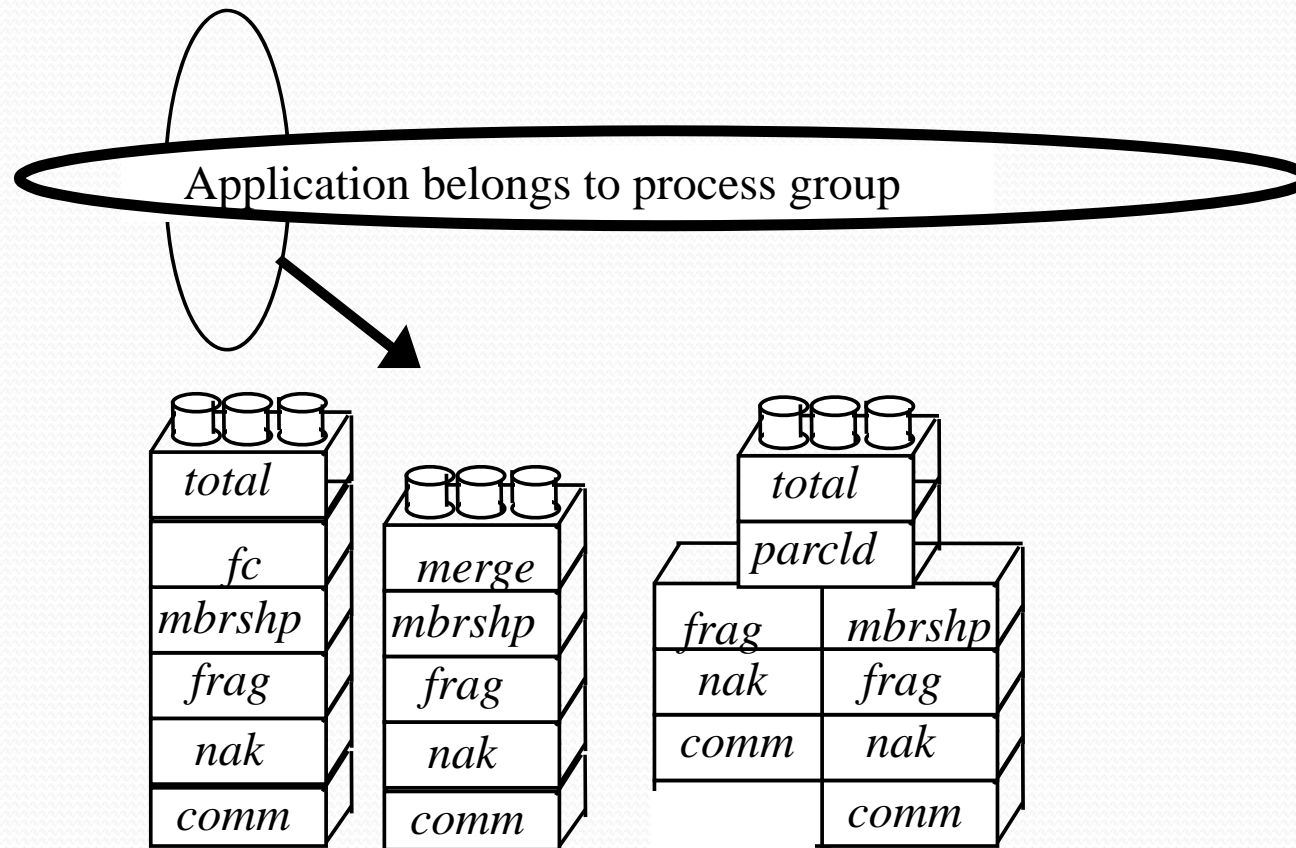


# Features of major virtual synchrony platforms

- Horus, JGroups and Ensemble
  - All were developed at Cornell: successors to Isis
  - These focus on flexible protocol stack linked directly into application address space
    - A stack is a pile of micro-protocols
    - Can assemble an optimized solution fitted to specific needs of the application by plugging together “properties this application requires”, lego-style
    - The system is optimized to reduce overheads of this compositional style of protocol stack
  - JGroups is very popular.
  - Ensemble is somewhat popular and supported by a user community. Horus works well but is not widely used.



# Horus/JGroups/Ensemble protocol stacks







# JGroups (part of JBoss)

- Developed by Bela Ban
  - Implements group multicast tools
    - Virtual synchrony was on their “to do” list
    - But they have group views, multicast, weaker forms of reliability
  - Impressive performance!
  - Very popular for Java community
- Downloads from [www.JGroups.org](http://www.JGroups.org)





# Spread Toolkit

- Developed at John Hopkins
  - Focused on a sort of “RISC” approach
    - Very simple architecture and system
    - Fairly fast, easy to use, rather popular
  - Supports one large group within which user sees many small “lightweight” subgroups that seem to be free-standing
  - Protocols implemented by Spread “agents” that relay messages to apps





# Quicksilver

- Under development here at Cornell
- Remarkably scalable, stable, fast
- But current emphasis is on use within Live Objects framework, not as a general purpose toolkit
- Focus of the research right now: a “properties framework” in which protocol properties
  - Are specified in a high-level language
  - Which compiles to a specialized runtime infrastructure





# What about Paxos?

- Protocol developed by Lamport
  - Initial version was for a fixed group membership, and allowed any process to propose an update
  - Later refined with a leader-driven version very similar to the GMS view update protocol
- Lamport proposed an elegant safety proof that remains a classic
- Actual Paxos implementations need to deal with dynamic membership and are much more complex, resemble virtual synchrony





# What about Paxos?

- Basically, Paxos offers the user a safe (dynamically uniform) abcast protocol
  - As you would expect, this is quite slow compared to asynchronous non-dynamically uniform fbcast...
  - ... But modern computers are *very fast* and perhaps the performance hit isn't such an issue
- Google uses this in Chubby, a lock service
- Microsoft also uses it, in SQL server clusters (also for locking). But they use virtual synchrony in the Enterprise Cluster Manager itself.





# Virtual Synchrony, Paxos... and *beyond*?

- Work underway at Microsoft Silicon Valley is seeking to unify the models
  - They are talking about “virtually synchronous Paxos”
  - Goal is to have
    - Dynamic membership with view synchrony
    - Paxos as the basic protocol for new events
    - Proof system, as in the case of basic Paxos, to let people prove protocols correct and reason about their applications





# Summary?

- Role of a toolkit is to package commonly used, popular functionality into simple API and programming model
- Group communication systems have been more popular when offered in toolkits
  - If groups are embedded into programming languages, we limit interoperability
  - If groups are used to transparently replicate deterministic objects, we're too inflexible
- Many modern systems let you match the protocol to your application's requirements