FutureGRI D: A Program for long term research into GRI D Systems Architecture

Jon Crowcroft, Steve Hand, Tim Harris, I an Pratt

The Computer Laboratory, University of Cambridge +

Andrew Herbert, Director, Microsoft Research Cambridge

O. Introduction

- r Program of work between the Computer Lab, and Microsoft Research
- r Builds on existing collaborationsr Designed as a set of loosely couple basic research
- projects r Common elements to projects, which lead to
- understanding r Later, full systems architecture will emerge for a Future GRI D.
- PhD studentships efficient use of funds (and to be honest, we have more good applicants than money ∠

1. Who, where, how, what

- r Collaborative tools based on Scribe and Pastry instead (or as well as) I P multicast (P2P CSCW) (existing RFC on PGM etc)
- r Search based on locality and on partial content matching (publications this month)
- Computation based on large scale systems and massively redundant partition of computational problems (a.k.a. spread spectrum)
- Extension of Pasta work on mutable, persistent P2P storage (publications)







Concerns with IP Multicast

Scalability with number of groups

r

- m Routers maintain per-group state
- ${\tt m}$ Analogous to per-flow state for QoS guarantees
- m Aggregation of multicast addresses is complicated
- r Supporting higher level functionality is difficult
 - m IP Multicast: best-effort multi-point delivery service
 - $\tt m$ End systems responsible for handling higher level functionality $\tt m$ Reliability and congestion control for IP Multicast complicated
- r Inter-domain routing is hard.
- r No management of flat address space.
- r Deployment is difficult and slow
 - m ISP's reluctant to turn on IP Multicast

End System P2P Multicast

Why is self-organization hard?

- Dynamic changes in group membership
- m Members join and leave dynamically
- m Members may die
- Limited knowledge of network conditions
- ${\tt m}$ Members do not know delay to each other when they join ${\tt m}$ Members probe each other to learn network related information
- m Overlay must self-improve as more information available
- Dynamic changes in network conditions
- m Delay between members may vary over time due to congestion Use Pastry/Scribe P2P system as it provides precisely these charactistics...







12

Vector Space Search: applications

- Resource discovery:
 - m Points represent resource requirements of jobs and resource availability of machines
 - m Nodes act as brokers between jobs and systems that can host them
- r Network position could be reflected in the broker's coordinates
 - m Promote scalability through disjoint operation of user communities when requests are satisfied by local facilities

13

15

Spread Spectrum Computing -Project 3

- r Use redundancy coding *ideas*
- r For code and data,
- r Dissemination uses high degrees of replication
- r Collection of responses is m Distributed (P2P)
- m Fault tolerant (like <u>SETL@Home</u> and the set of ideas in a lot of cryptanalysis work recently r Highly Optimised Tolerance (c.f. John

14

16

Doyle's work at CalTech).

Global Storage - Project 4

- r Available anywhere, anytime and fast!
- r Must cope with node and network failures m Use replication, information dispersal codes
- r Must cope with `flash crowds' m Automatic load balancing and distribution
- r Must allow local caching for performance
- m Challenge of maintaining consistency r Must provide `hands free' administration
 - m Self-organizing system

Global Storage with Pasta

- r Uses P2P Distributed Hash Table techniques m More complex structures necessary? B*trees?
- r Aims to provide traditional file-system like semantics (incl. efficient mutability, quotas)
- r Also, wider look at shared workspaces to support ad-hoc collaboration
 - m Not all participants fully trusted...
 - m Need versioning, `views' and 'overlaying'
 - m Object-specific locking and atomicity enforced by storage system