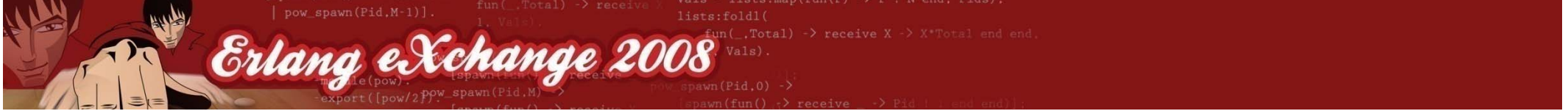


Building a transactional distributed data store with Erlang

Alexander Reinefeld, Florian Schintke, Thorsten Schütt

Zuse Institute Berlin,
onScale solutions GmbH



Transactional data store - What for?

- **Web 2.0 services:** shopping, banking, gaming, ...
 - don't need full SQL semantics, key/value DB often suffice
 - e.g. Mike Stonebreaker: "One size does *not* fit all"

- **Scalability matters**
 - $>10^4$ accesses per sec.
 - many concurrent writes

Slashdot NEWS FOR NERDS. STUFF THAT MATTERS.

► [Log In](#) | [Create Account](#) | [Help](#) | [Subscribe](#) | [Firehose](#)

▼ **Sections**

- [Main](#)
- [Apple](#)
- [AskSlashdot](#)
- [Backslash](#)
- [Books](#)
- [Developers](#)
- [Games](#)
- [Hardware](#)
- [Interviews](#)
- [IT](#)
- [Linux](#)

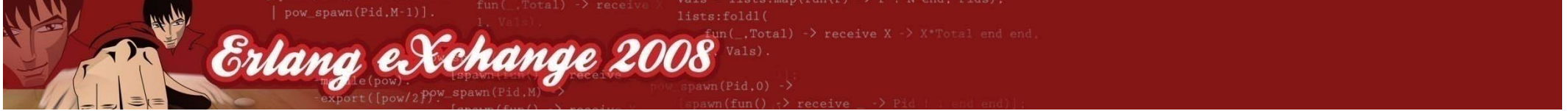
Is the One-Size-Fits-All Database Dead?

Posted by [kdawson](#) on Tue Jan 09, 2007 10:50 PM
from the [specialized-and-optimized](#) dept.

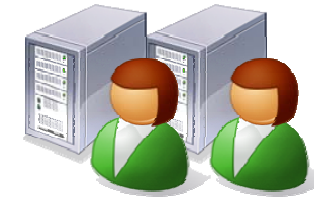
[jlbrown](#) writes

"In a new benchmarking paper, MIT professor Mike Stonebraker and colleagues demonstrate that [specialized databases can have dramatic performance advantages over traditional databases](#) (PDF) in four areas: text processing, data warehousing, stream processing, and scientific and intelligence applications. The advantage can be a factor of 10 or higher. The paper includes some interesting 'apples to apples' performance comparisons between commercial implementations of specialized architectures and relational databases: data warehousing and stream processing."





Traditional Web 2.0 hosting

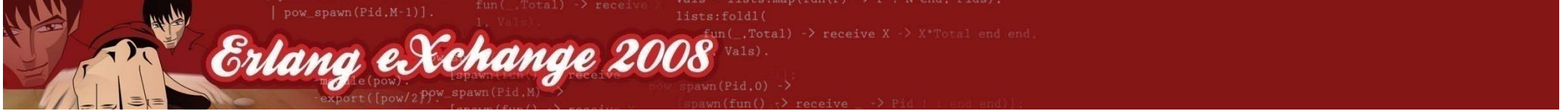


Clients

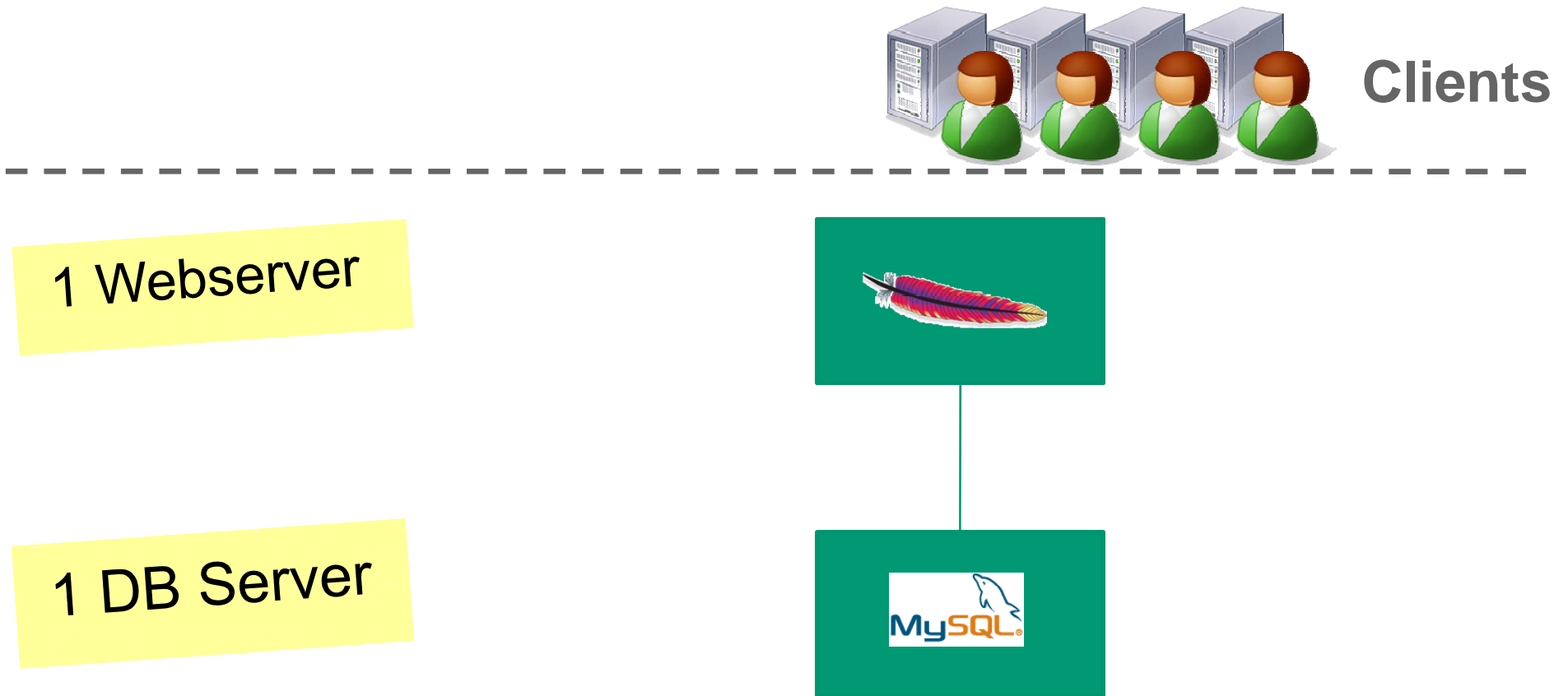
1 Server

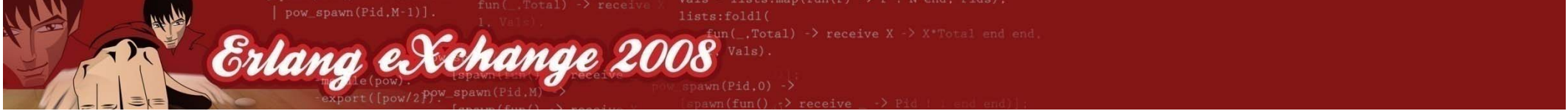
- Webserver
- DB Server



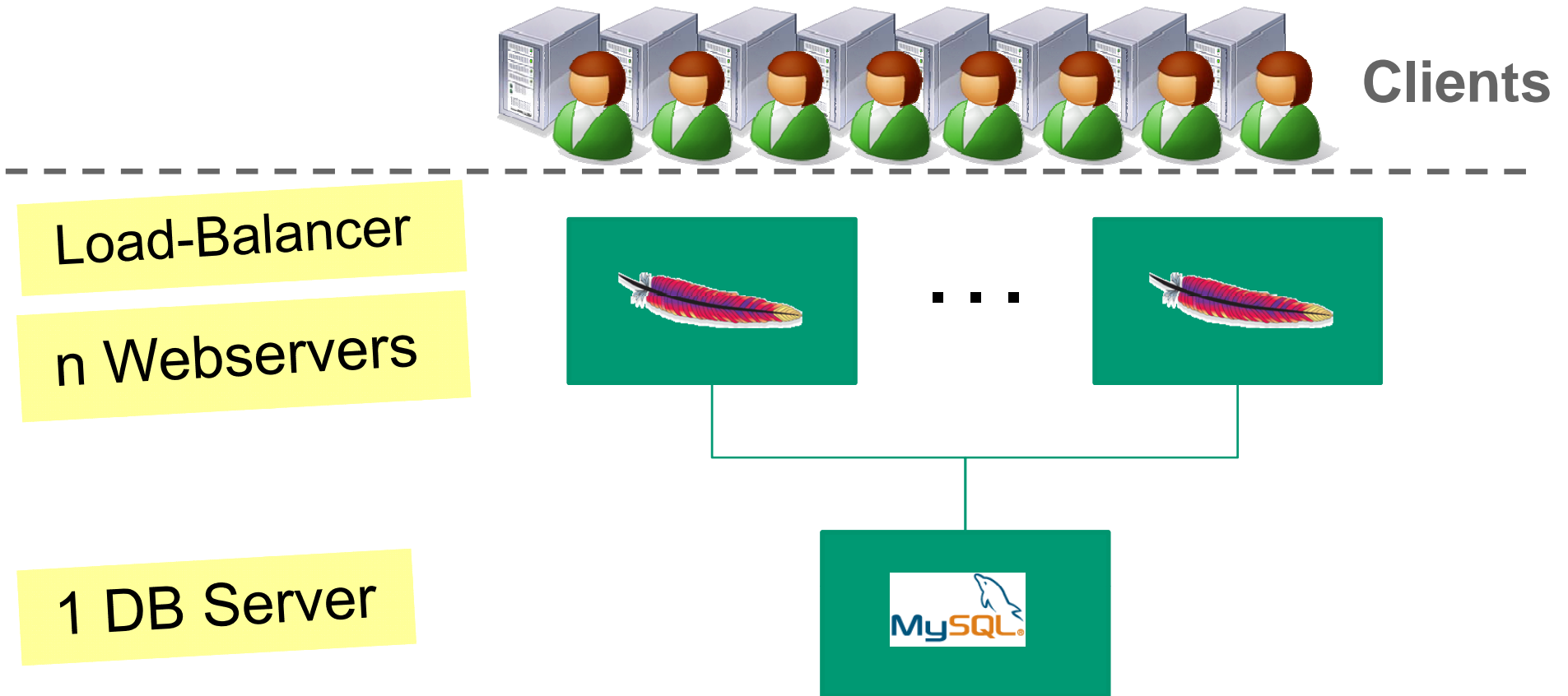


Traditional Web 2.0 hosting

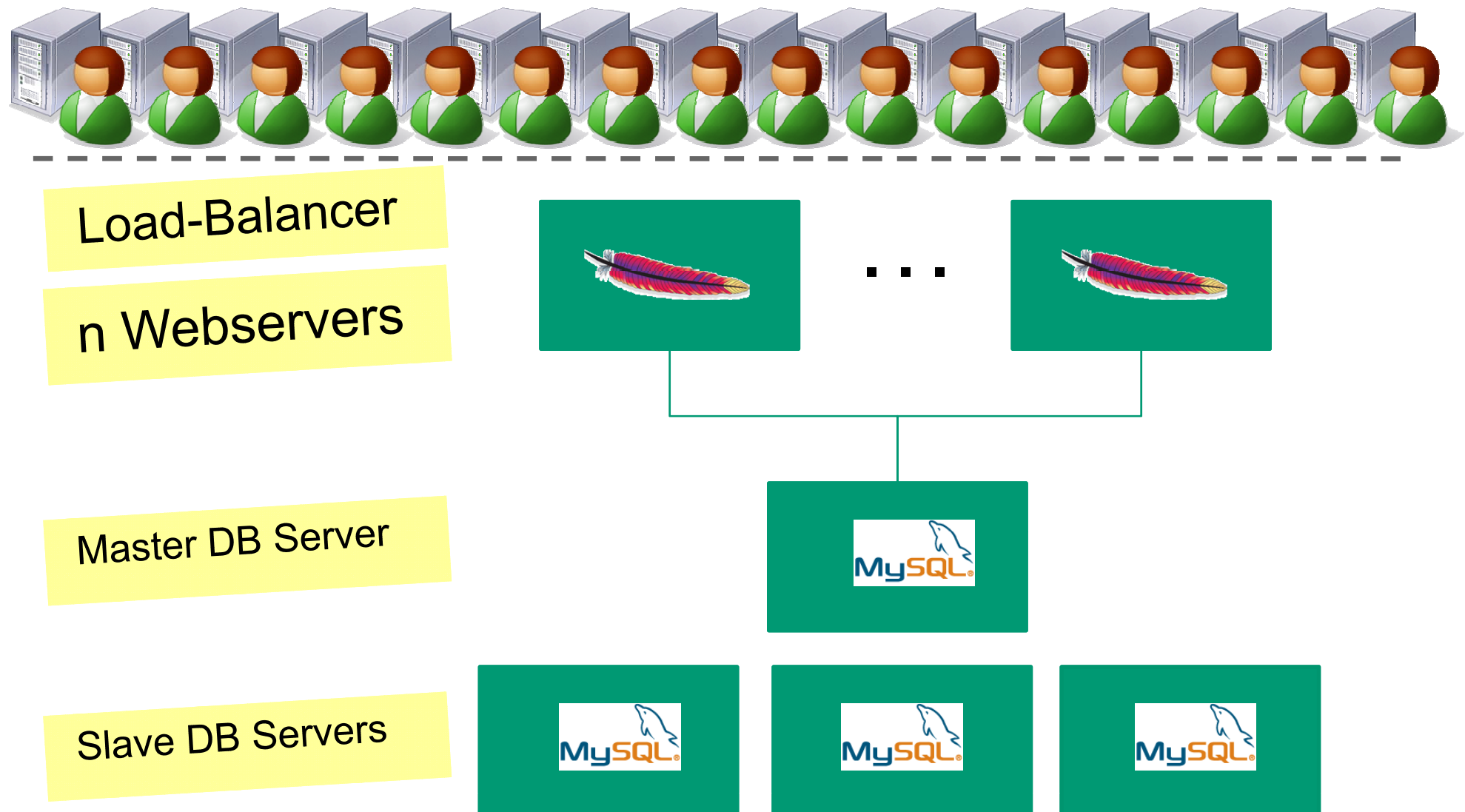


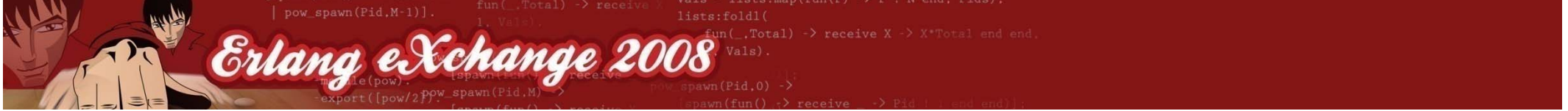


Traditional Web 2.0 hosting



Traditional Web 2.0 hosting



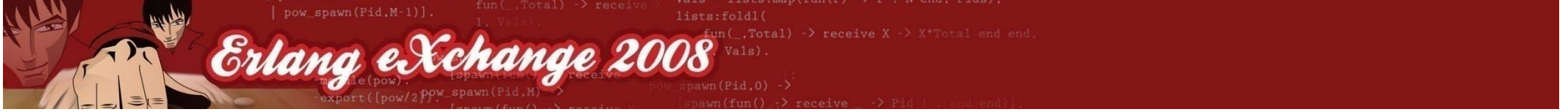


Now think big. Really BIG.

Not how fast our code is today, but:

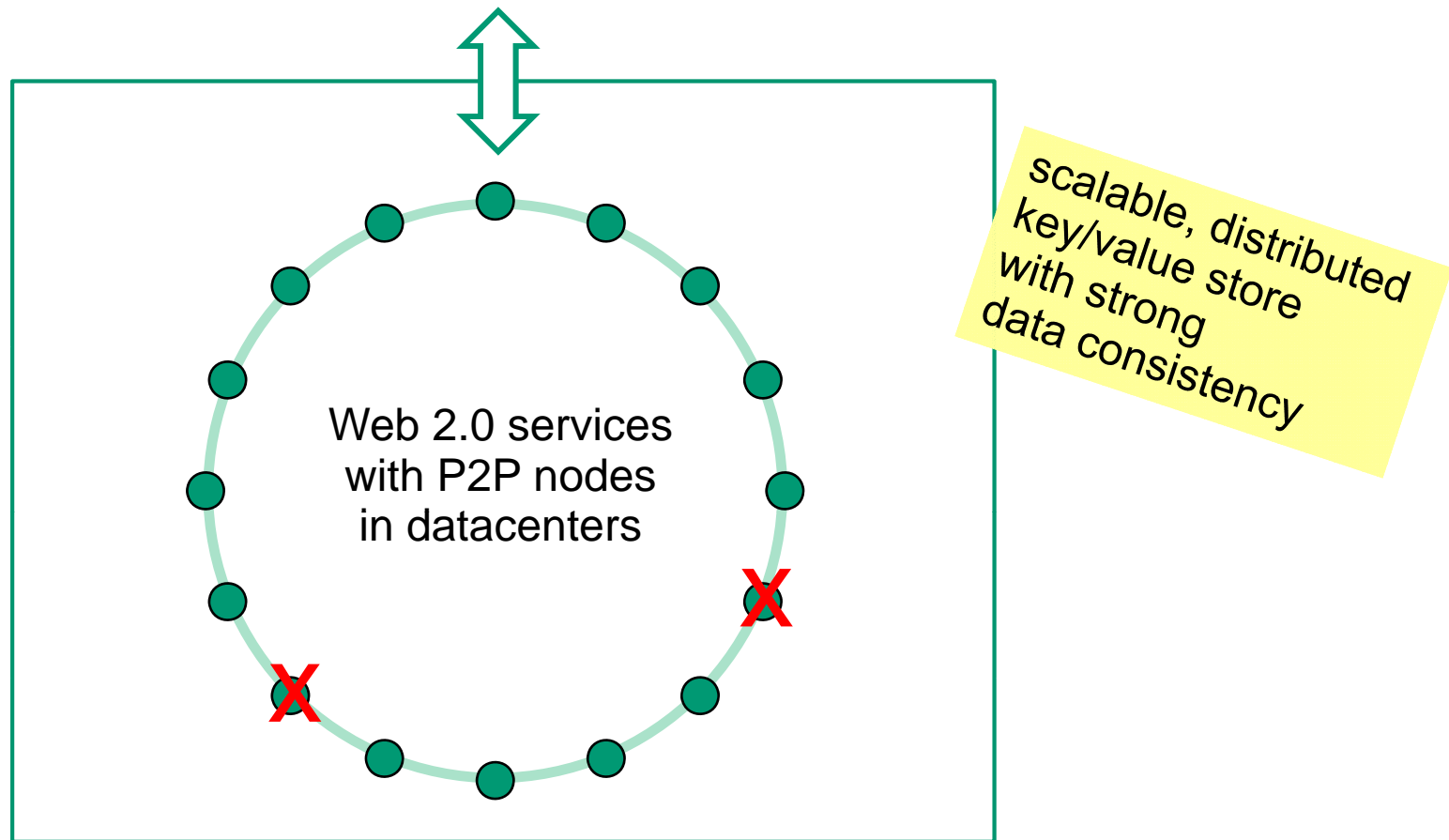
- Can it “scale out”?
- Can it run in parallel? ... distributed?
- Any common resources causing locking?

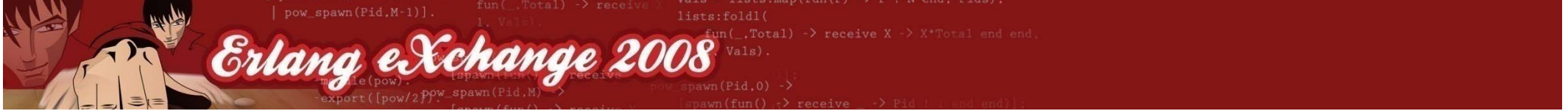
Asymptotic performance matters!



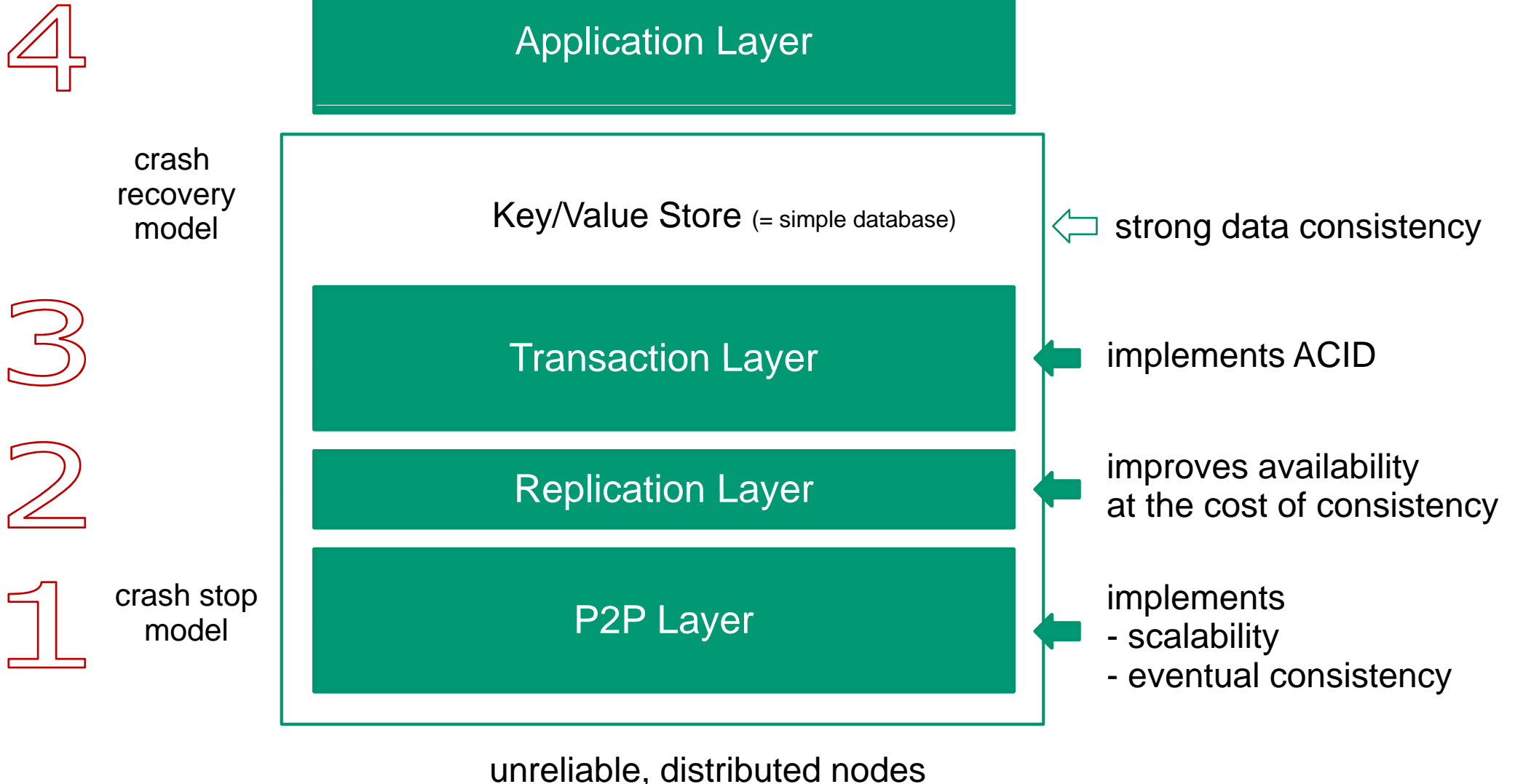
Our Approach: P2P makes it scalable

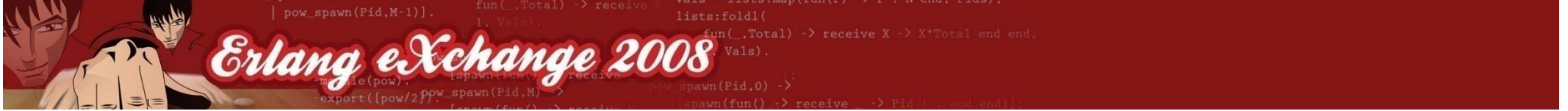
“arbitrary” number of clients





Our Approach

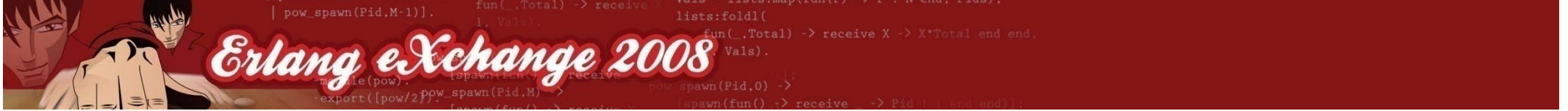




1

providing a scalable distributed data store:

P2P LAYER

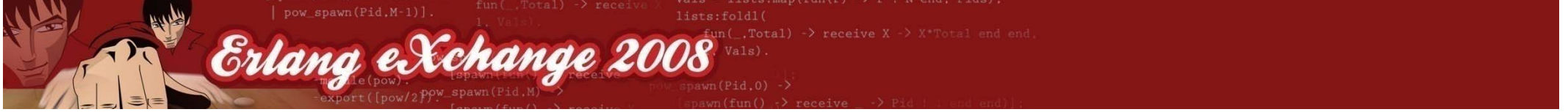


Key/Value Store

- for storing “items” (= “key/value pairs”)
 - synonyms: “key/value store”, “dictionary”, “map”, ...
- just 3 ops
 - insert(key, value)
 - delete(key)
 - lookup(key)

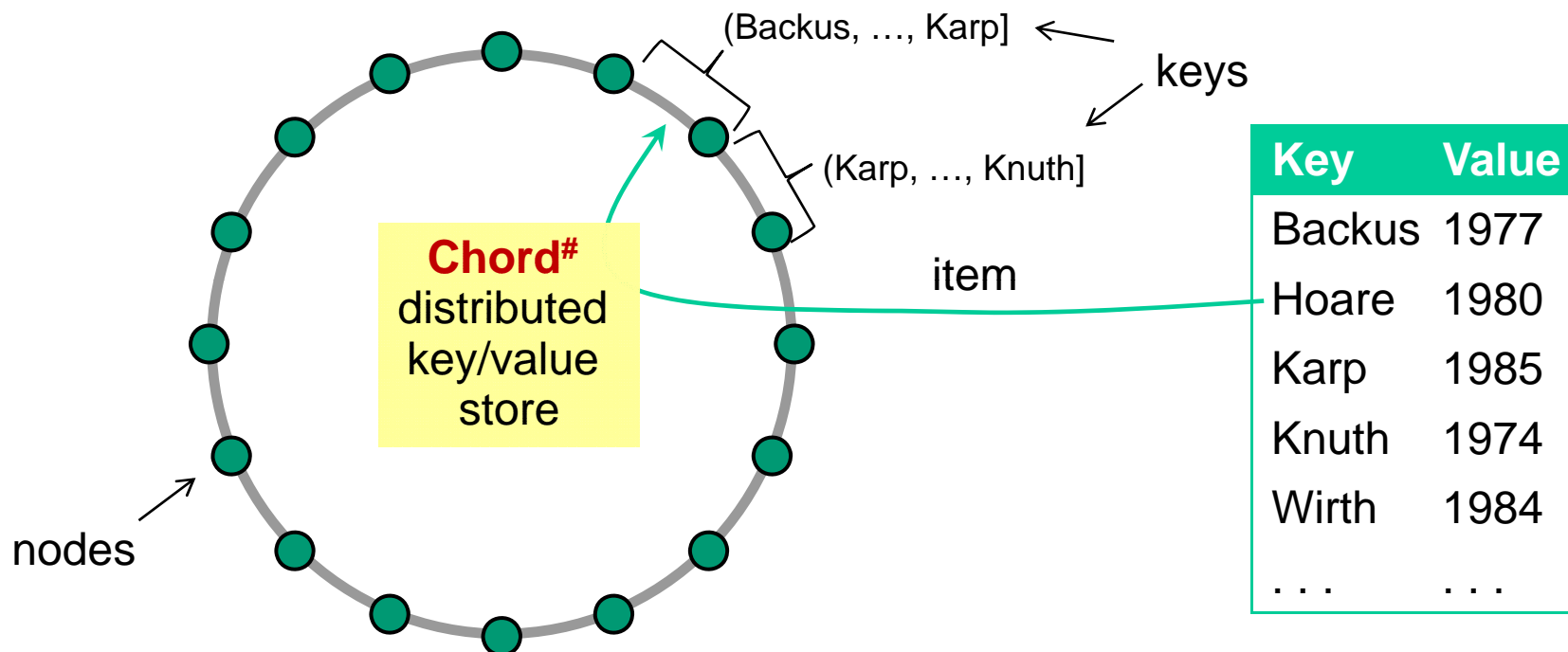
Turing Award Winners

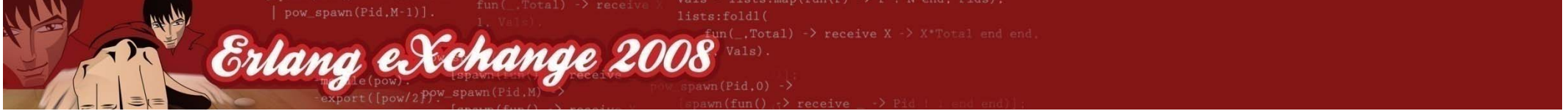
Key	Value
Backus	1977
Hoare	1980
Karp	1985
Knuth	1974
Wirth	1984
...	...



Chord# - Distributed Key/Value Store

- **key space:** total order on items (strings, numbers, ...)
- nodes have a random key as their position in ring
- items are stored on the successor node (clockwise)



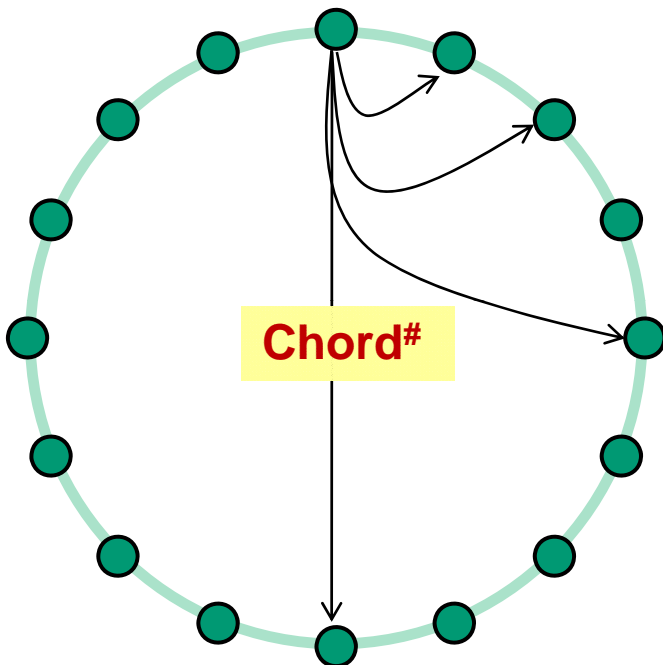


Routing Table and Data Lookup

Building the routing table

- $\log_2 N$ pointers
- exponentially spaced pointers

$$pointer_i = \begin{cases} successor & : i = 0 \\ pointer_{i-1} . pointer_{i-1} & : i \neq 0 \end{cases}$$

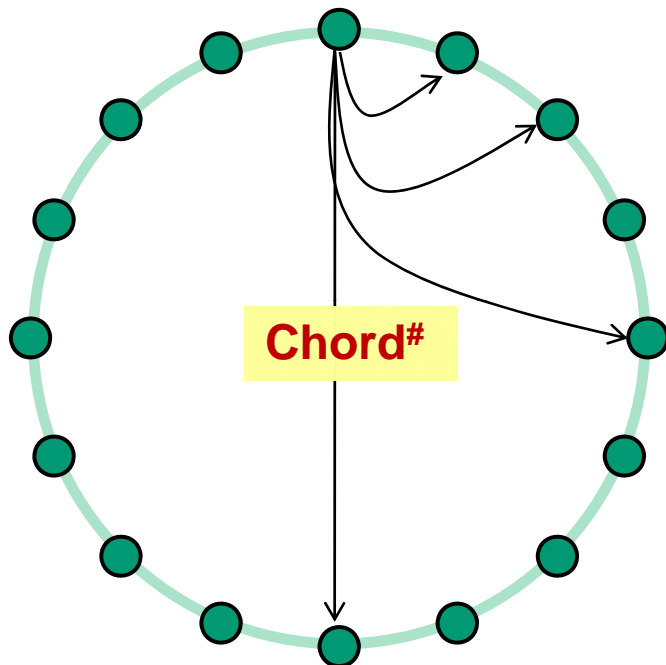


Routing Table and Data Lookup

Building the routing table

- $\log_2 N$ pointers
- exponentially spaced pointers

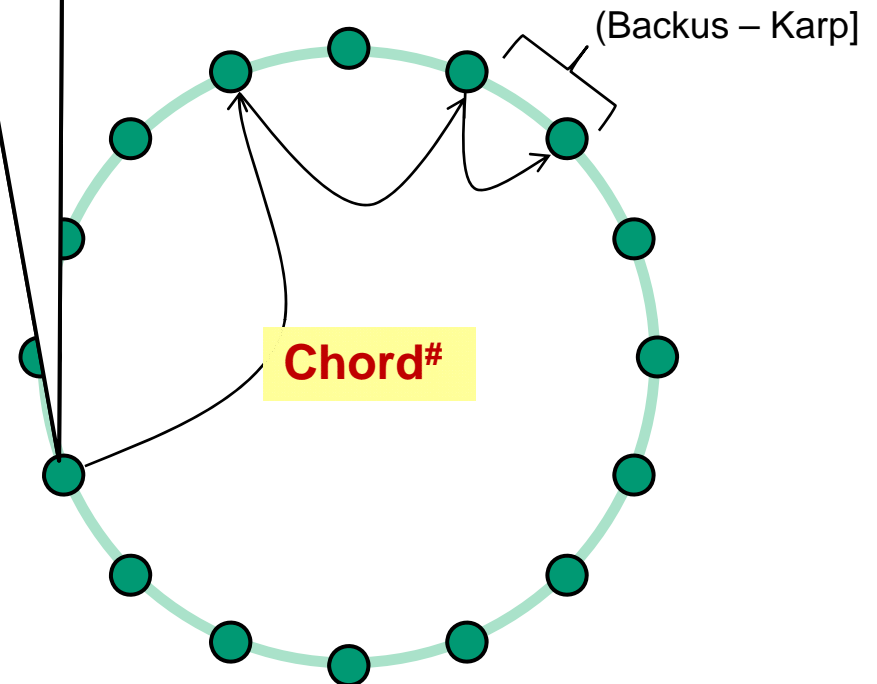
$$pointer_i = \begin{cases} successor & : i = 0 \\ pointer_{i-1} . pointer_{i-1} & : i \neq 0 \end{cases}$$

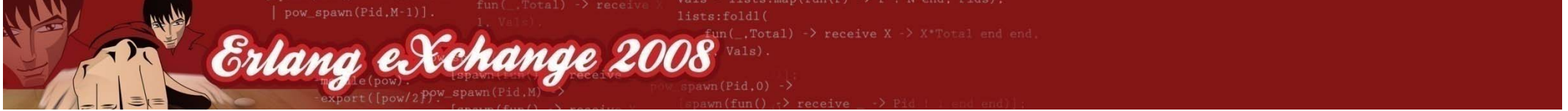


Retrieving items

- $\leq \log_2 N$ hops

Example:
lookup (Hoare)
started from here



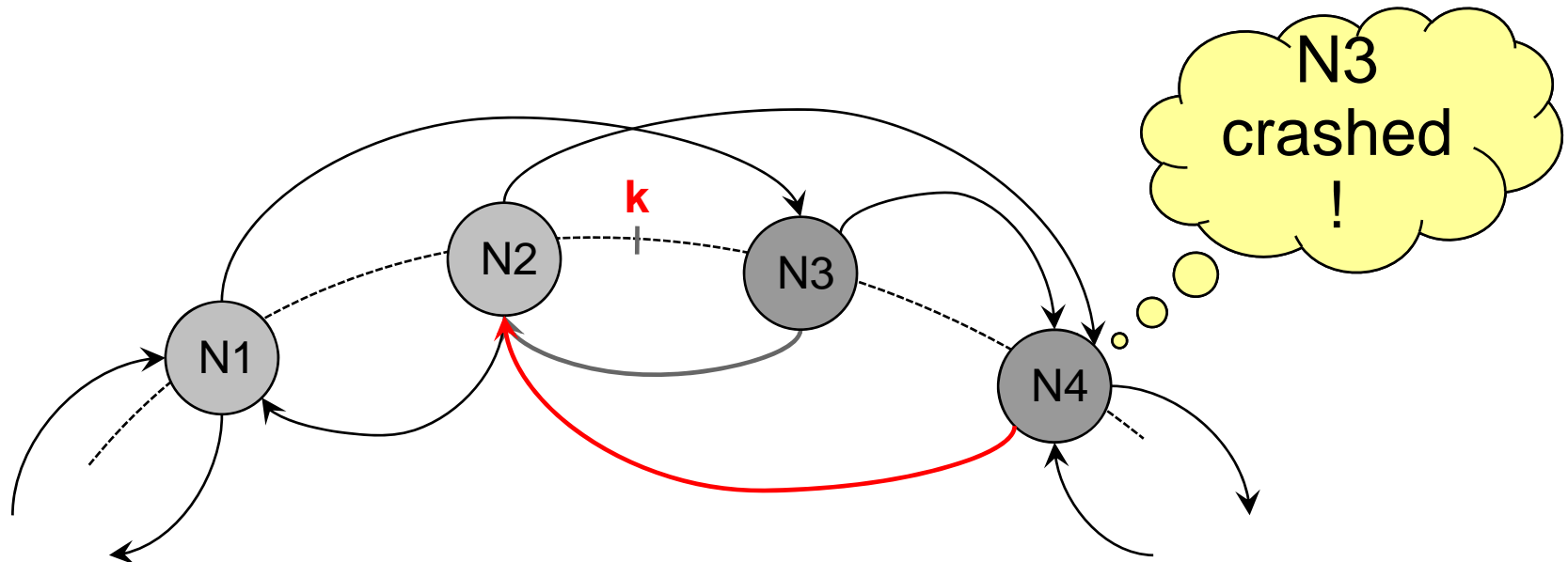


Churn

- Nodes *join, leave, or crash* at any time
- Need “**failure detector**” to check aliveness of nodes
 - failure detector may be wrong: Node dead? Or just slow network?
- Churn may cause inconsistencies
 - need local repair mechanism

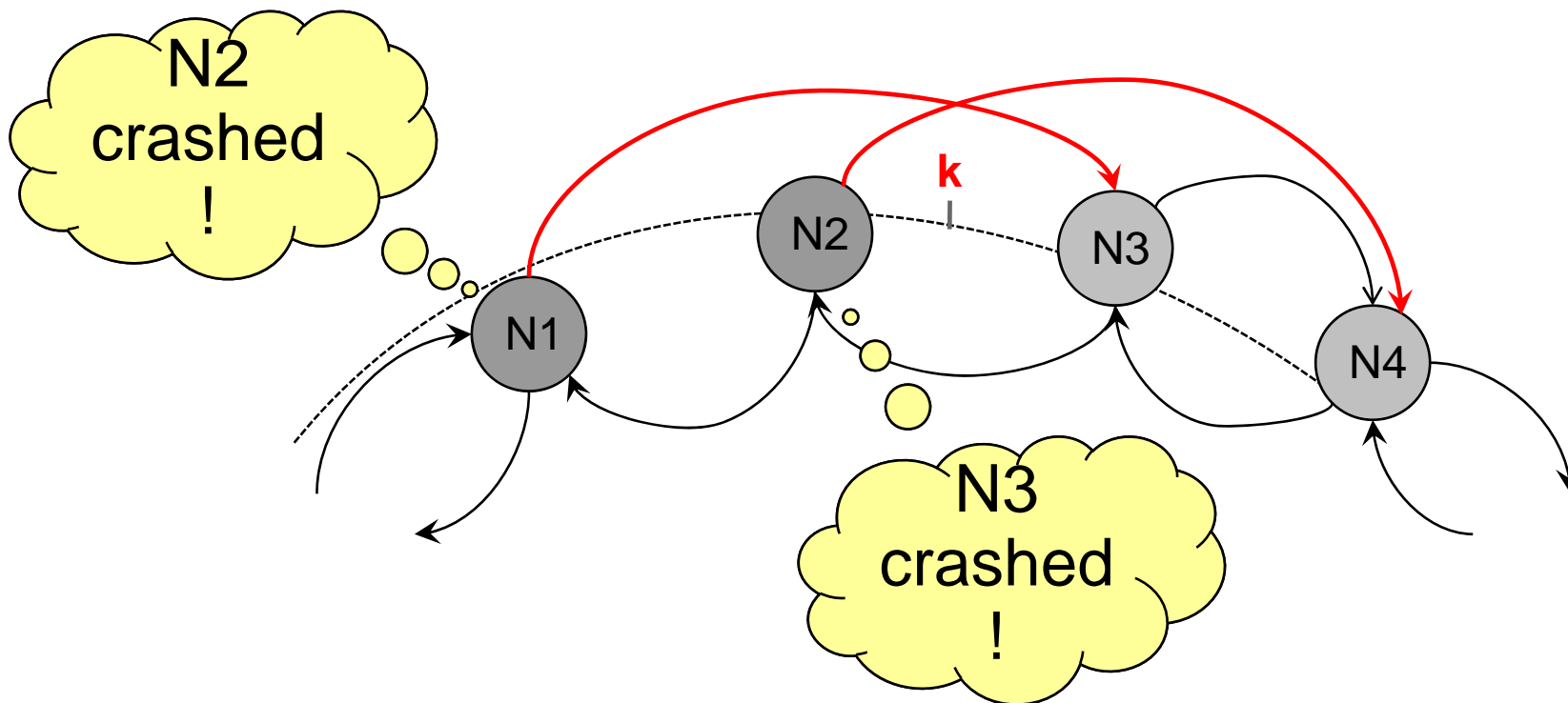
Responsibility Consistency

- **Violated responsibility consistency** caused by imperfect failure detector: Both, N3 *and* N4 claim responsibility for item **k**



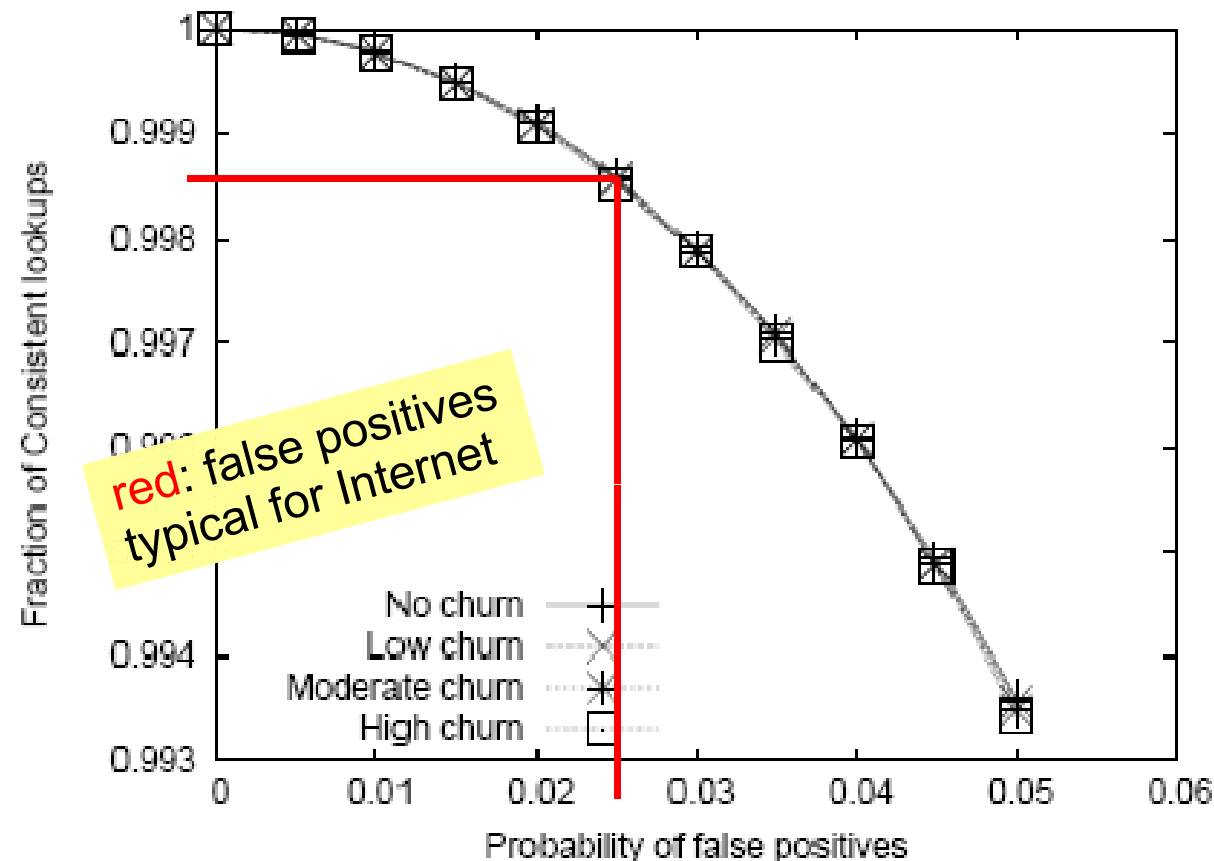
Lookup Consistency

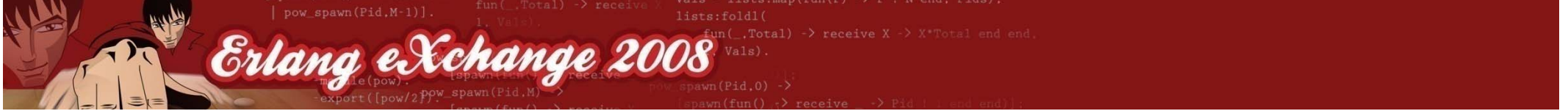
- **Violated lookup consistency** caused by imperfect failure detector:
lookup(k): at N1 → N3, but at N2 → N4



How often does this occur?

- Simulated nodes with imperfect failure detectors
(A node detects another alive node as dead probabilistically)



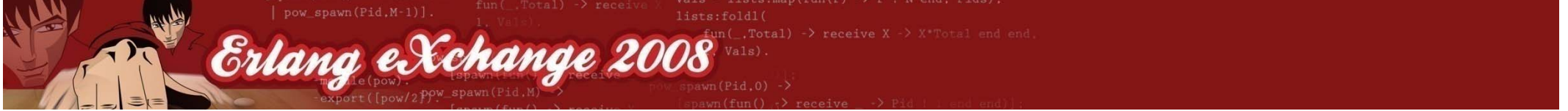


SUMMARY

1

P2P LAYER

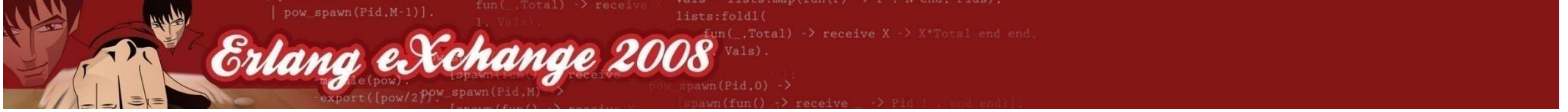
- Chord# provides a key/value store
 - scalable
 - efficient: $\log_2 N$ hops
- Quality of failure detector is crucial
- Need **replication** to prevent data loss ...



2

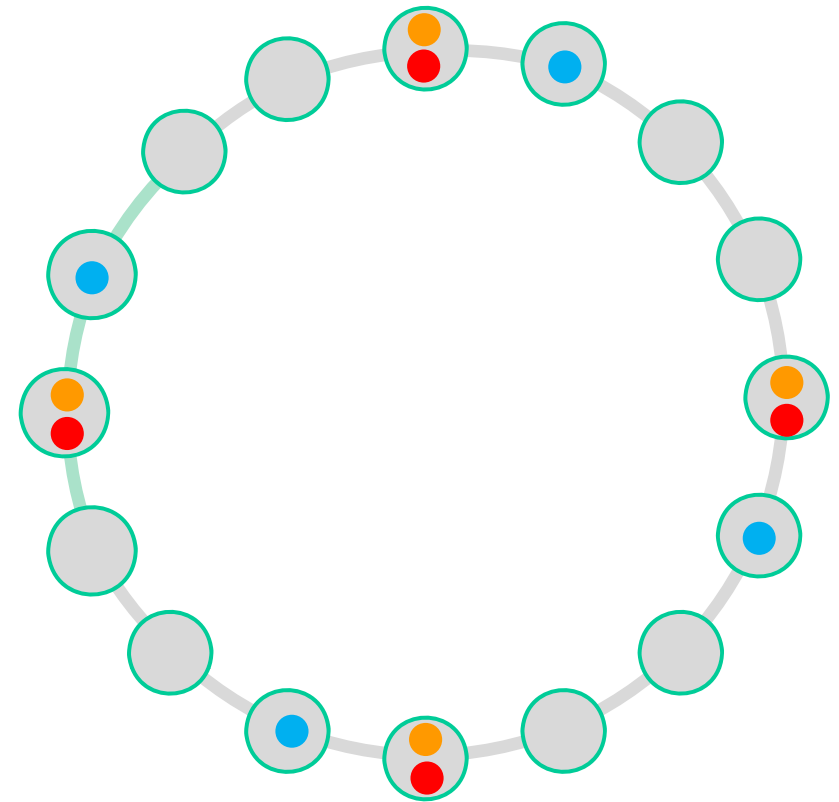
improving availability

REPLICATION LAYER



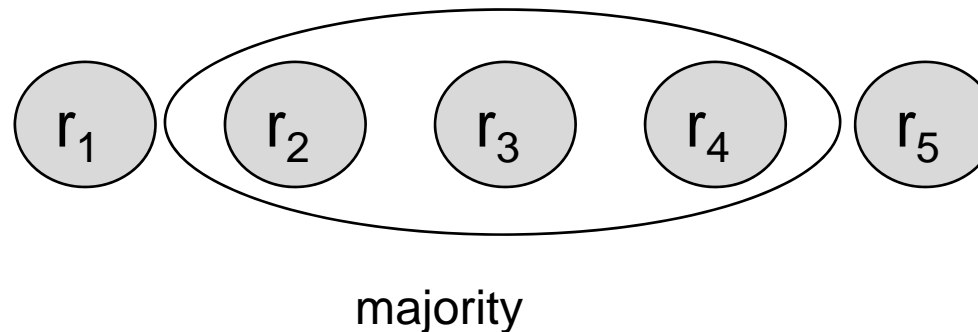
Replication

- Many schemes
 - symmetric replication →
 - succ. list replication
 - ...
- Must ensure data consistency
 - need quorum-based methods

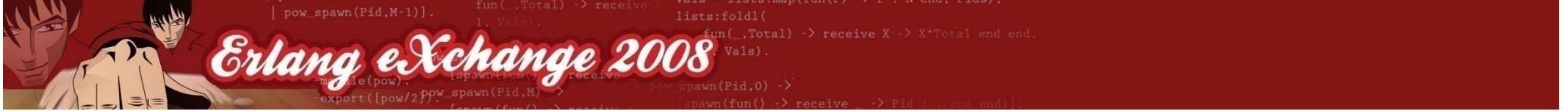


Quorum based algorithms

- Enforce consistency by operating on majorities



- Comes at the cost of increased latency
 - but latency can be avoided by clever replica distribution in datacenters (**cloud computing**)

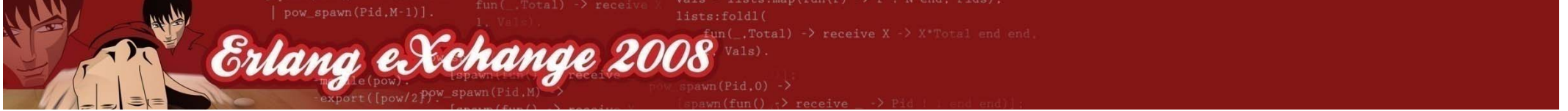


SUMMARY

2

REPLICATION LAYER

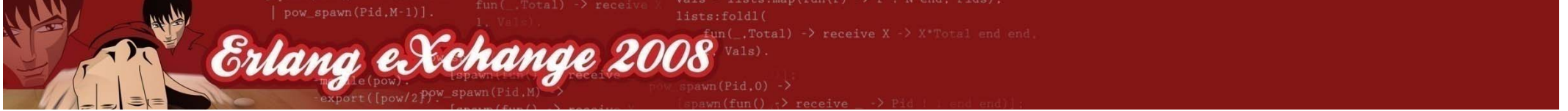
- availability in face of churn
- quorum algorithms
- But need **transactional** data access ...



3

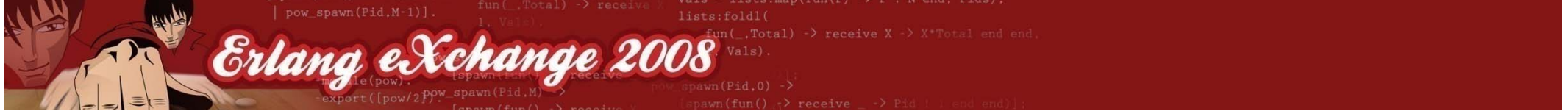
coping with concurrency:

TRANSACTION LAYER



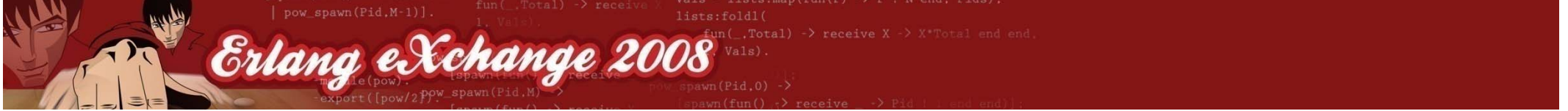
Transaction Layer

- Transactions on P2P are challenging because of ...
 - **churn**
 - changing node responsibilities
 - **crash stop** fault model
 - as opposed to crash recovery in traditional DBMS
 - imperfect **failure detector**
 - don't know whether node crashed or slow network



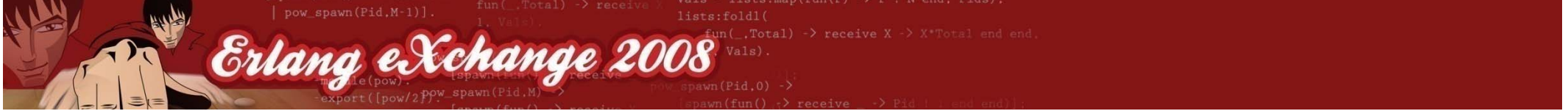
Strong Data Consistency

- What is it?
 - When a write is finished, all following reads return the new value.
- How to implement?
 - Always read/write **majority** $\lfloor f/2 \rfloor + 1$ of f replicas.
 - Latest version is always in the read or write set
 - Must ensure that replication degree is $\leq f$



Atomicity

- What is it?
 - Make **all** or **no** changes!
 - Either ‘commit’ or ‘abort’.
- How to implement?
 - 2PC? Blocks if the transaction manager fails.
 - 3PC? Too much latency.
 - We use a variant of the **Paxos Commit** Protocol
 - non-blocking: Votes of transaction participants are sent to multiple “acceptors”

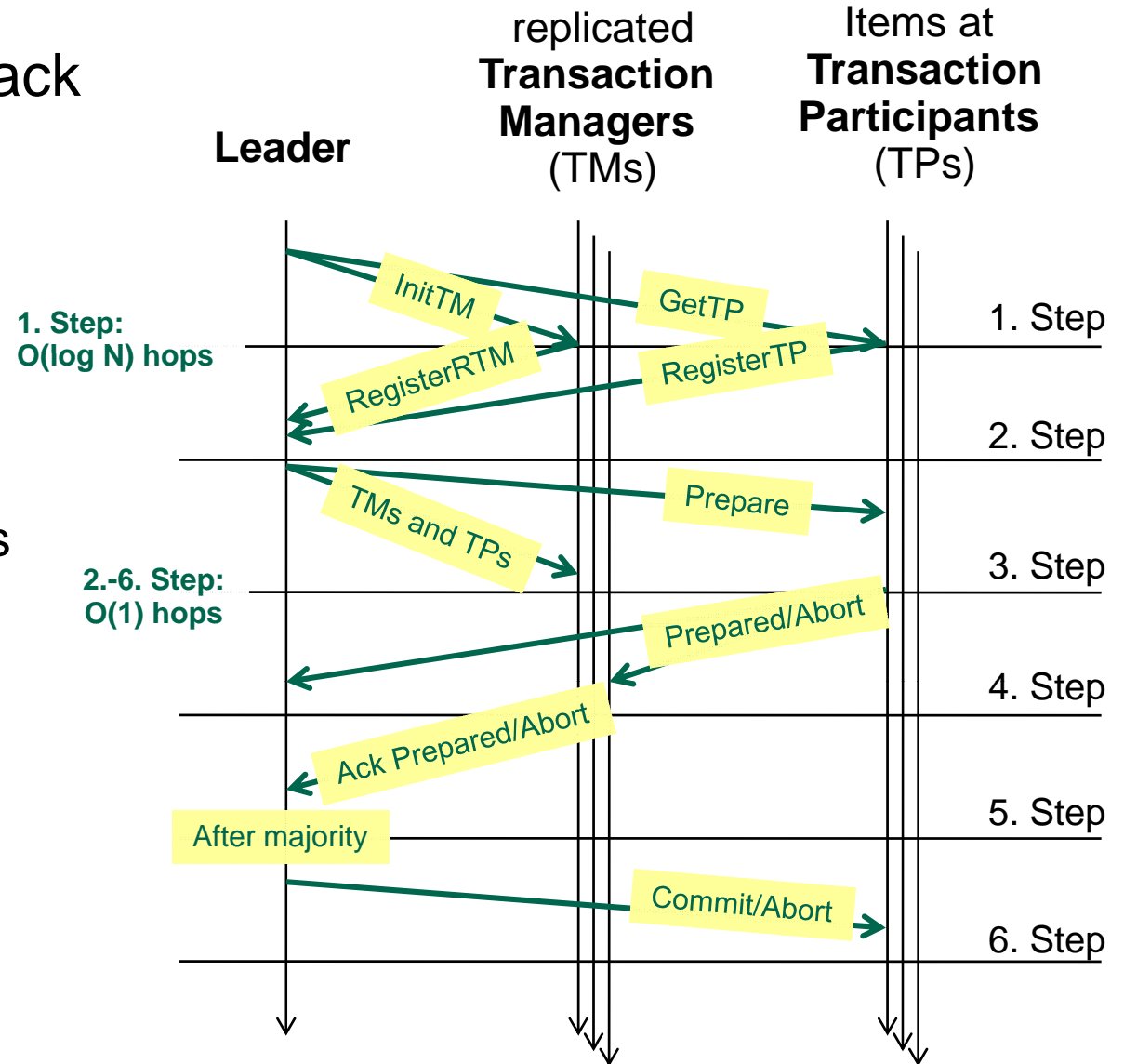


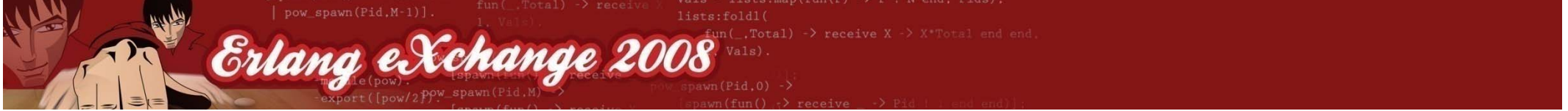
Adapted Paxos Commit

- Optimistic CC with fallback
- Write
 - 3 rounds
 - non-blocking (fallback)
- Read even faster
 - reads majority of replicas
 - just 1 round
- succeeds when $>f/2$ nodes alive

Adapted Paxos Commit

- Optimistic CC with fallback
- Write
 - 3 rounds
 - non-blocking (fallback)
- Read even faster
 - reads majority of replicas
 - just 1 round
- succeeds when $>f/2$ nodes alive





Transactions have two purposes:

Consistency of replicas & consistency across items

User Request

BOT

- debit (a, 100);
- deposit (b, 100);

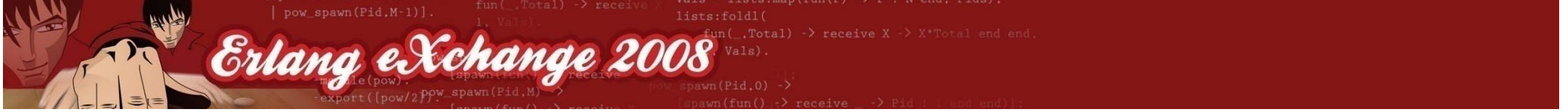
EOT

Operation on replicas

BOT

- debit (a1, 100);
- debit (a2, 100);
- debit (a3, 100);
- deposit (b1, 100);
- deposit (b2, 100);
- deposit (b3, 100);

EOT

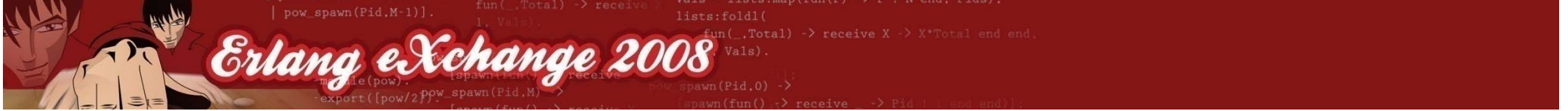


SUMMARY

3

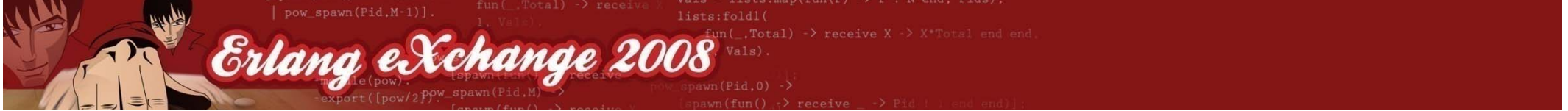
TRANSACTION LAYER

- Consistent update of items and replicas
- Mitigates some of the overlay oddities
 - node failures
 - asynchronicity



4

demonstrator application:
WIKIPEDIA



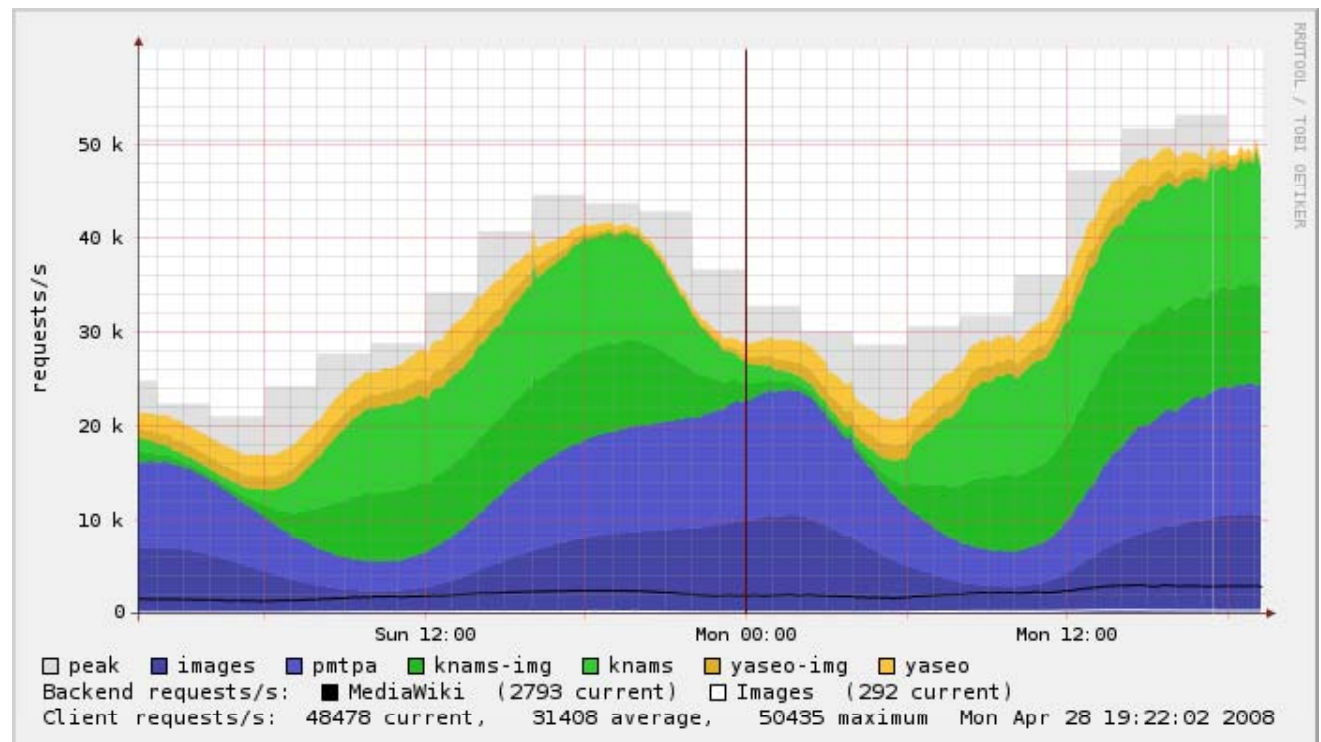
Wikipedia

Top 10 Web sites

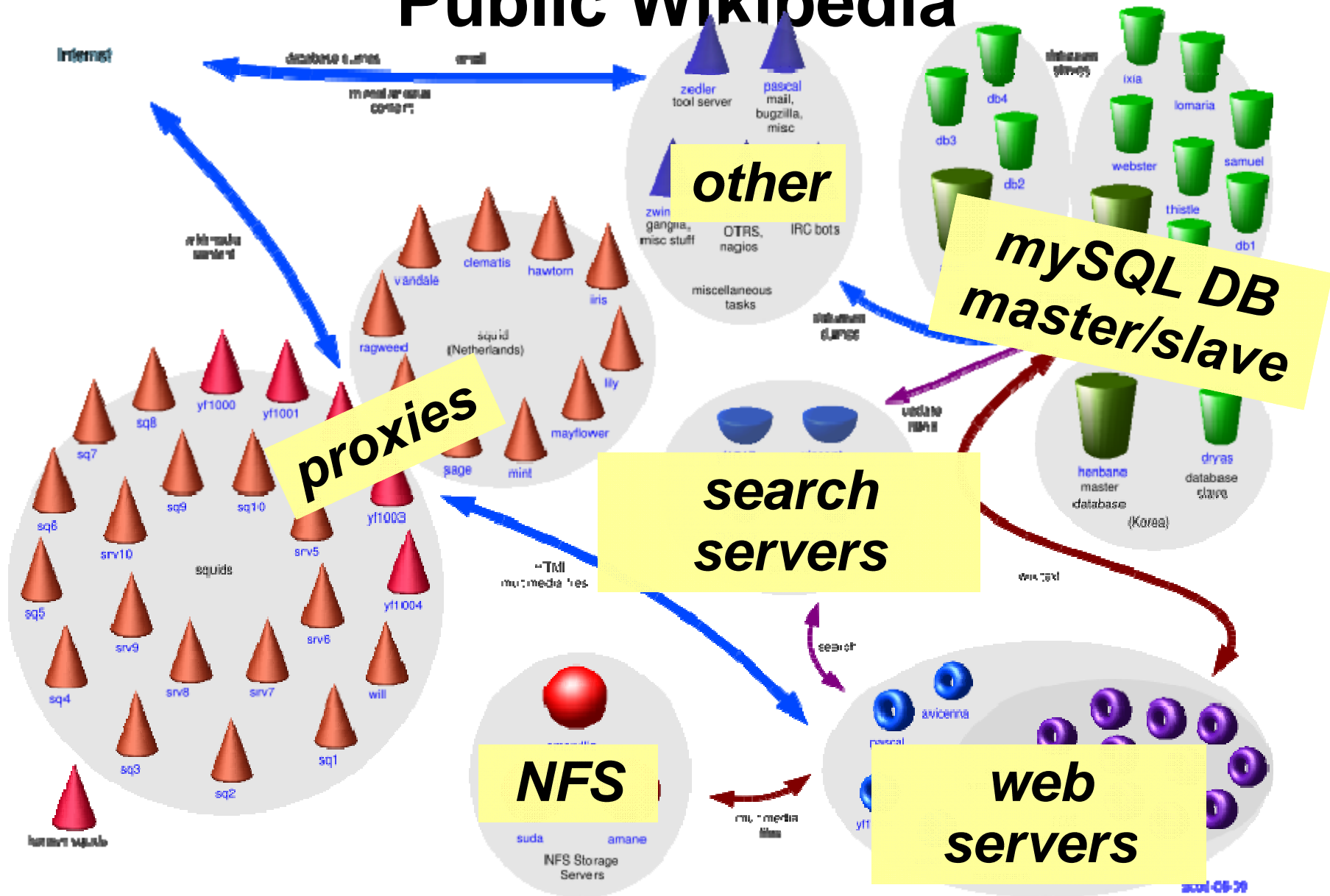
1. Yahoo!
2. Google
3. YouTube
4. Windows Live
5. MSN
6. Myspace
7. **Wikipedia**
8. Facebook
9. Blogger.com
10. Yahoo!カテゴリ

50.000 requests/sec

- 95% are answered by squid proxies
- **only 2,000 req./sec hit the backend**



Public Wikipedia



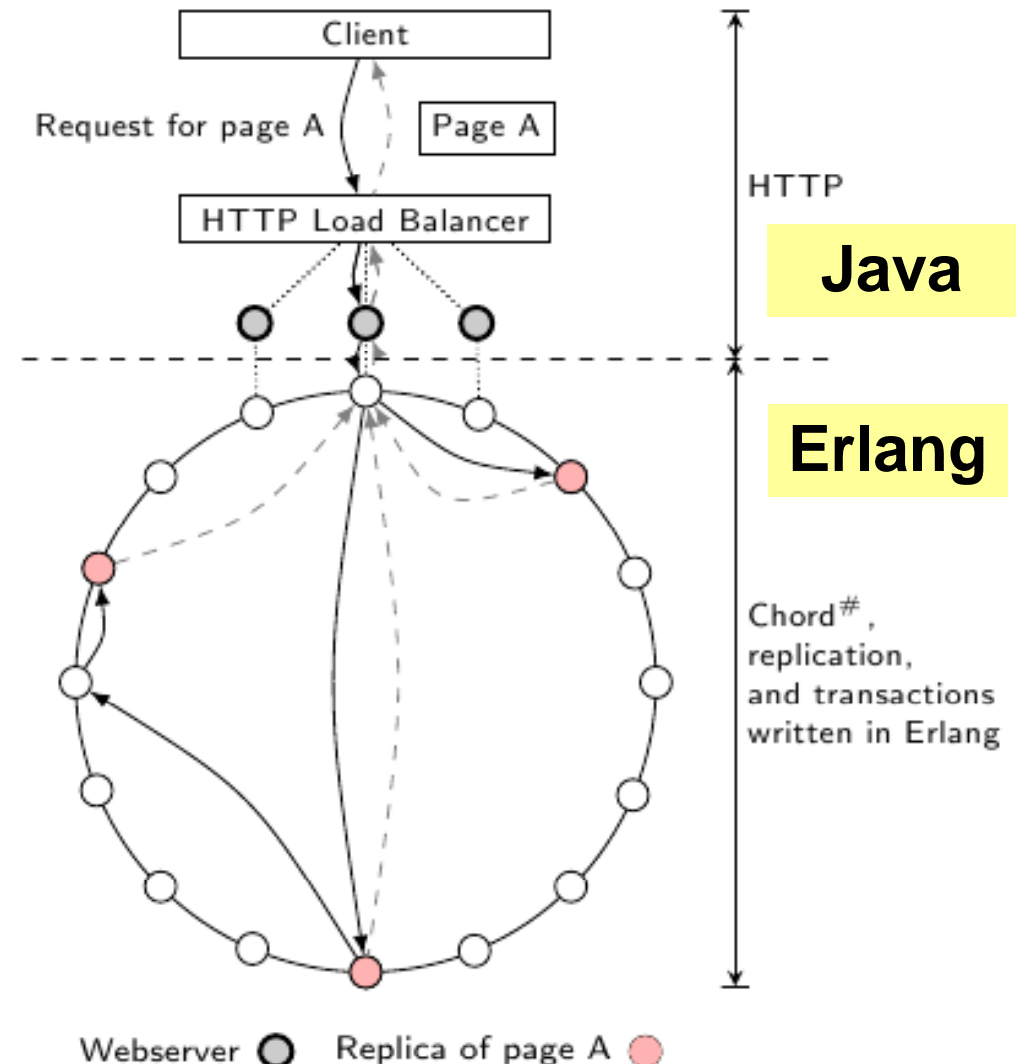
Our Wikipedia

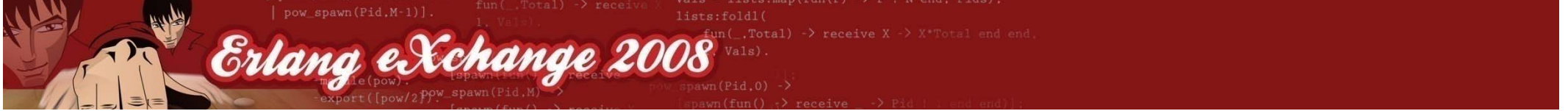
Renderer

- Java
 - Tomcat, Plog4u
- Jinterface
 - Interface to Erlang

Key/Value Store

- Chord[#] + Replication + Transactions





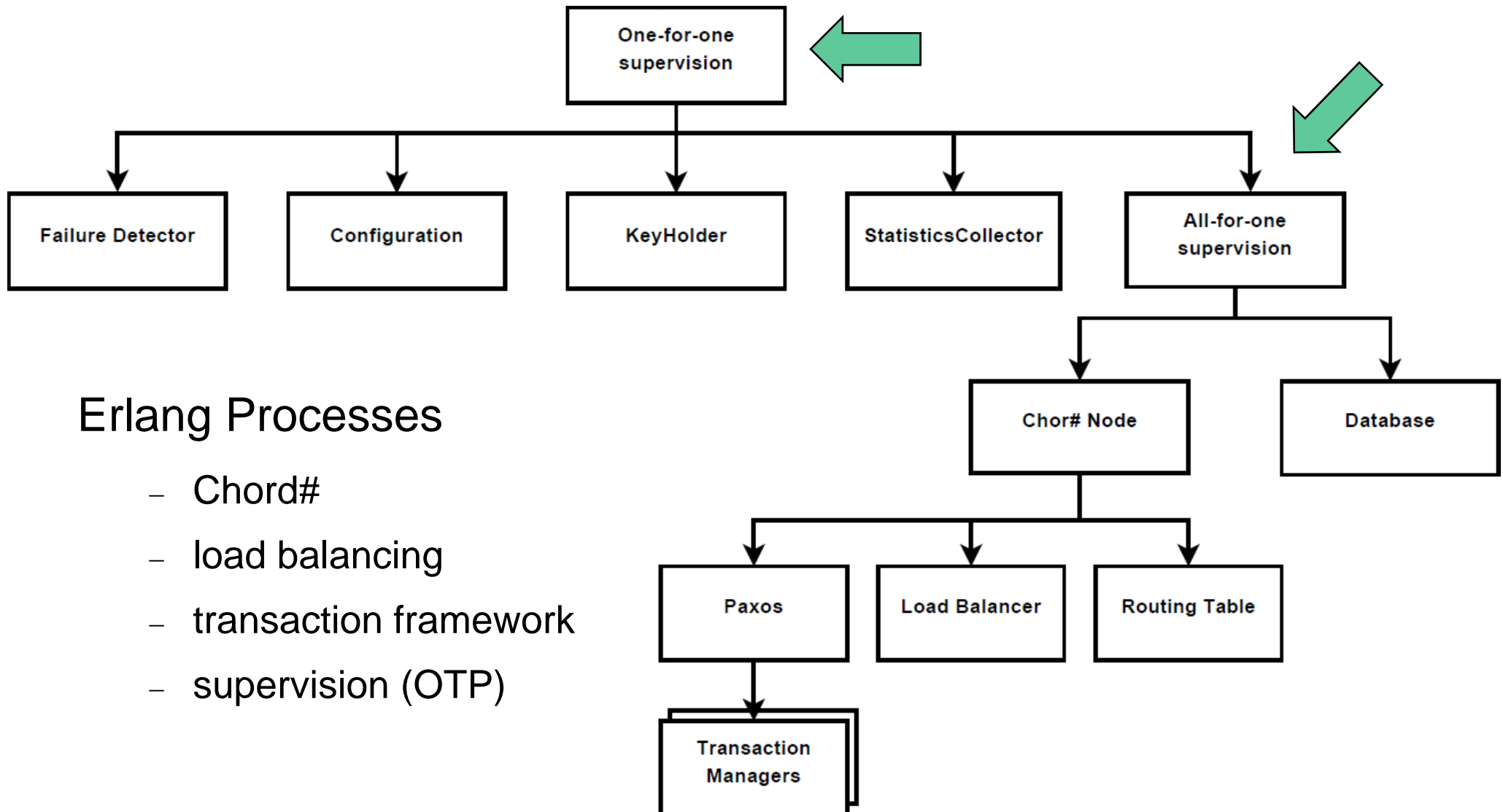
Mapping Wikipedia to Key/Value Store

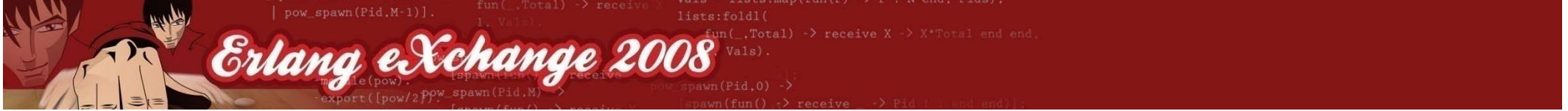
- Mapping

	key	value
page content	title	list of Wikitext for all versions
backlinks	title	list of titles
categories	category name	list of titles

- For each insert or modify we must
 - update backlinks
 - update category page(s)
- } *write transaction*

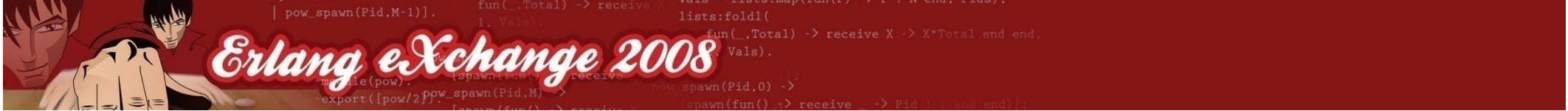
Erlang Processes





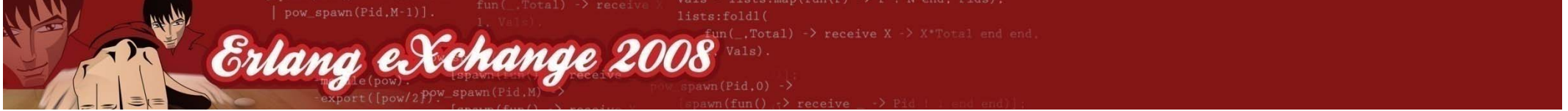
Erlang Processes (per node)

- **Failure Detector** supervises Chord# nodes and sends crash messages when a failure is detected.
- **Configuration** provides access to the configuration file and maintains parameter changes made at runtime.
- **Key Holder** stores the identifier of the node in the overlay.
- **Statistics Collector** collects statistics information and forwards them to statistic servers.
- **Chord# Node** performs the main functionality of the node, e.g. successor list and routing table.
- **Database** stores key/value pairs in each node.



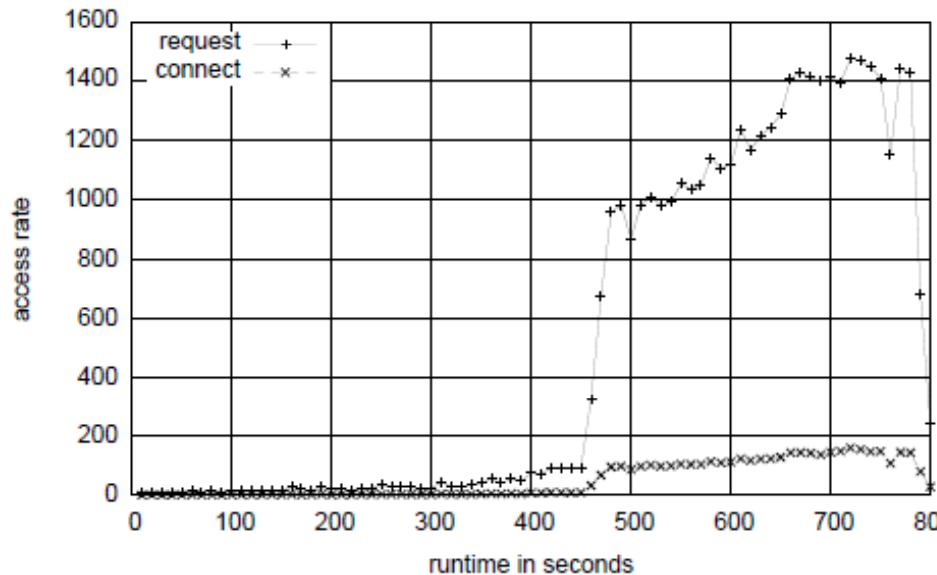
Accessing Erlang Transactions from Java via Jinterface

```
void updatePage(string title, int oldVersion, string newText)
{
    Transaction t = new Transaction();           //new transaction
    Page p = t.read(title);                       // read old version
    if (p.currentVersion != oldVersion)           // concurrent update?
        t.abort();
    else {
        t.write(p.add(newText));                  // write new text
        //update categories
        foreach(Category c in p)
            t.write(t.read(c.name).add(title));
        t.commit();
    }
}
```

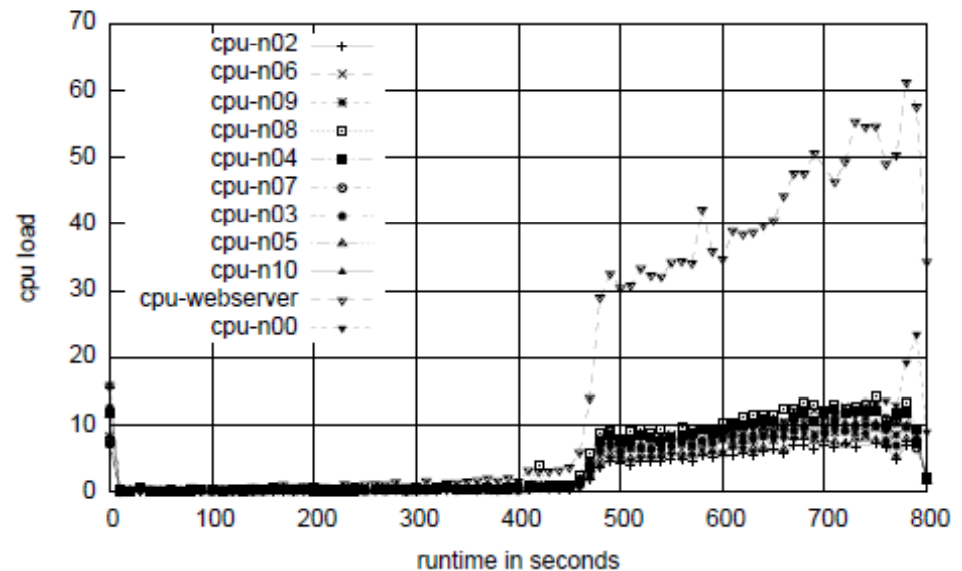


Performance on Linux Cluster

test results with load generator



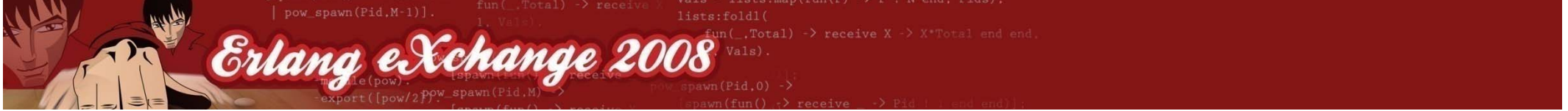
throughput with increasing access rate over time



CPU load with increasing access rate over time

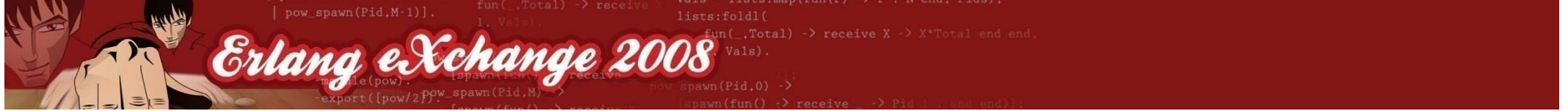
1500 trans./sec on 10 CPUs

2500 trans./sec on 16 CPUs (64 cores) and 128 DHT nodes



Implementation

- 11,000 lines of Erlang code
 - 2,700 for transactions
 - 1,300 for Wikipedia
 - 7,000 for Chord# and infrastructure
- Distributed Erlang
 - currently has weak security and limited scalability
 - ⇒ we implemented own transport layer on top of TCP
- Java for rendering and user interface



Hauptseitrn - Wikipedia - Konqueror


Location Edit View Go Bookmarks Tools Settings Window Help

Location: W http://localhost.zib.de:8080/wiki?title=Hauptseitrn

Google Search

javac ... Yaws Chord#... Login ... Login ... Overvi... W Error ... UPGRAD... First ... STDLIB... Chord#... W Haupts...

Oomeidn



WIKIPEDIA

De freie Enzyklopädie

navigation

- Hauptseitrn
- Zuefalls-Artikl
- Inhoitsvazeichnis

mitarbeit

- Wikipedia-Postal
- Fehlende Artikl
- Letzte Änderungen

hilf

- Fragen?
- Stammtisch
- Finanzielle Hilf

suach

Artikl Such

weakzeigkistn

- Links auf de Seitrn
- Valinks prüfn
- Nauflädn
- Spezial-Seitrn
- Seitrn zum Ausdrucken
- Bschtländige URL
- Seitrn zitirn

renderer

- Standard
- Koa Renderer
- Erlang Ausdruck


artikl bschprecha werklñ versionen/autoren

Hauptseitrn

Griß Gott in da boärischn Wikipedia! D' Wikipedia is a Projekt fürn Aufbau vò äna frein Enzyklopädie, in mehr wiä 200 Sprächñ. Dé Version is in da Boärischn Spräch gschriebm. Älle, dé an Dialekt redn, der dā dazuäghert (in Äidbayan, Östareich und Südtiroi), derfm mitschreibm. Egäl, obs auf Nord-, Mittl- oda Südboärisch is. Äis, wås d' schreibst, derf **frei kopiärt und weitagebm** wern. Da **Äfang** is gānz äfäch!


- Eine kurze **Beschreibung** dieses Projekts in anderen Sprachen.
- A short **description** about this project in other languages.

{{NUMBEROFARTICLES}} Artikl: **A bis Z** und **Neie Ärtikl**




Geogràfie

Afrika - Amerika - Asien - Australien - Europa: Bayan • Deitschländ • Fränkreich • Großbritannien • Östareich • Isländ • Italien • Europäische Union • Slowénien • Liechtñstā • Belgien • Südtiroi




Gschicht

Vur- und Friägschicht • Äitatum • Mittlāita • Friärare Neizeit • Imperialismus und Wöidkriäge • Zwoäta Wöidkriäg




Gsöischäft

Politik • Wirtschaft • Recht • Ethik




Technik

Vakehr • Architektur • Elektrotechnik • Computer • Fotografie




Wissnschäft

Geisteswissnschäft • Philosophie • Thèologie • Ästronomie • Biologie • Kemie • Mathematik • Fysik • Ethnologie • Pädagogik • Sozialwissnschäft • Rechtswissnschäft




Glaubm

Buddhismus • Hinduismus • Judntum • Kristntum • Islām • Mythologie • Esoterik • Neiche Religionen



Sport

Eishockey • Fußbāi • Hāndbāi • Leichtätletik • Ténis • Biärkästnlaufm (Biärkastlāffa) • Schāfkopf • Krāxlñ

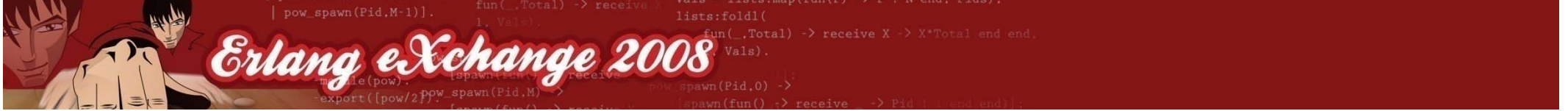


Kunst und Kultur

Fuim und Fernsehñ • Theata • Buidende Kunst • Buidhauarei • Musi • Kabarett • Biācha • Brauchtum • Spräch • Rockmusi • Gwānd • Musi vom Boärischn Sprächraum

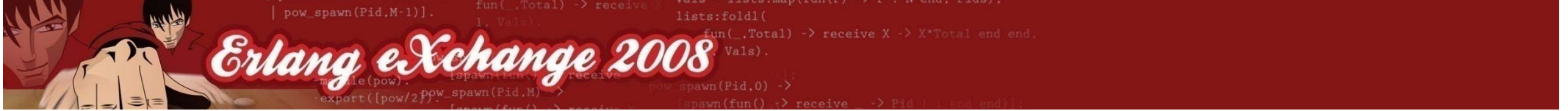
af: als: an: ast: be: bg: br: ca: cdo: cs: cu: cy: da: pdc: de: el: en: eo: es: et: eu: fi: fr: fy: gl: ht: hu: id: ik: io: is: it: ja: ka: ko: ksh: kw: la: li: lt: lv: mk: ms: nap: nds: nl: nn: no: oc: pl: pt: rm: ro: ru: ru-sib: scn: sh: simple: sk: sl: sq: sr: sv: th: tl: tr: uk: za: zh: zh-classical: zh-min-nan: zh-yue:

Kategorie: Wikipedia



5

SUMMARY



Summary

- P2P as new paradigm for Web 2.0 hosting
 - we support consistent, distributed write operations.
- Numerous applications:
 - Internet databases, transactional online-services, ...



All Things Distributed

Werner Vogels' weblog on building scalable and robust distributed systems.

Eventually Consistent

By Werner Vogels on December 19, 2007 10:03 PM | [Permalink](#) | [Comments \(16\)](#) | [TrackBacks \(6\)](#)

Recently there has been a lot of discussion about the concept of *eventual consistency* in the context of data replication. In this posting I would like to try to collect some of the principles and abstractions related to large scale data replication and the trade-offs between high-availability and data consistency. I consider this work-in-progress as I don't expect to get every definition crisp the first time.

There are two ways of looking at consistency. One is from the developer / client point of view; how they observe data updates. The second way is from the server side; how updates flow through the system and what guarantees systems can give with respect to updates.

Historical

In an ideal world there would only be one consistency model; when an update is made all observers will see that update. The first time this surfaced as difficult to achieve was in the database systems in the late seventies. The best "period piece" on this topic is by [Bruce Lindsay](#), et al, "Notes on Distributed Databases", Research Report RJ2571 (33471), IBM Research, July 1979. The fundamental principles for database replication are laid



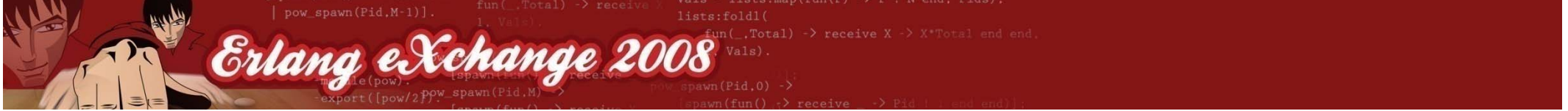
Contact Info

Werner Vogels

CTO - Amazon.com

werner@allthingsdistributed.com

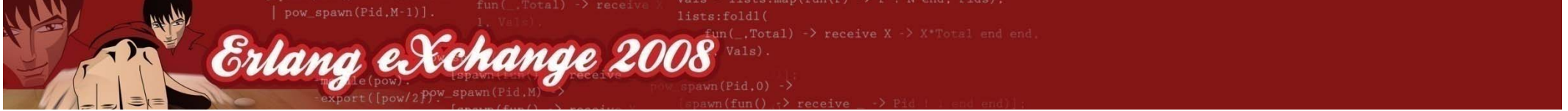
Tradeoff: High availability vs. data consistency



Team

- Thorsten Schütt
- Florian Schintke
- Monika Moser
- Stefan Plantikow
- Alexander Reinefeld
- Nico Kruber
- Christian von Prollius
- Seif Haridi (SICS)
- Ali Ghodsi (SICS)
- Tallat Shafaat (SICS)





Publications

Chord#

T. Schütt, F. Schintke, A. Reinefeld.

A Structured Overlay for Multi-dimensional Range Queries. Euro-Par, August 2007.

T. Schütt, F. Schintke, A. Reinefeld.

Structured Overlay without Consistent Hashing: Empirical Results. GP2PC, May 2006.

Transactions

M. Moser, S. Haridi.

Atomic Commitment in Transactional DHTs. 1st CoreGRID Symposium, August 2007.

T. Shafaat, M. Moser, A. Ghodsi, S. Haridi, T. Schütt, A. Reinefeld. **Key-Based Consistency and Availability in Structured Overlay Networks.** Infoscale, June 2008.

Wiki

S. Plantikow, A. Reinefeld, F. Schintke.

Transactions for Distributed Wikis on Structured Overlays. DSOM, October 2007.

Talks / Demos

IEEE Scale Challenge, May 2008
1st price (live demo)

