



Database Distribution

and why would I care about this as a developer..

Dirk Düllmann, IT-ADC

Database Workshop for LHC developers

24 January, 2005

Outline



Distributed Database Applications

- Why distributing database data?
- Which database technology can help to get to a distributed database application?
- Which distributed services are currently planned in LCG?
- How are applications development and deployment affected?



Distributed Deployment of Databases

- LCG provides infrastructure for storage, distribution and replication of file based data
- Physics applications (and grid m/w) require a similar services for data hosted in relational databases
 - Several applications and grid services use RDBMS - and more are coming this year
 - Several sites have already experience in providing RDBMS services
- Goals for the 3D project as part of LCG
 - provide LCG apps with a consistent, location independent **access to database services**
 - **increase the availability and scalability** of database applications via replication
 - arrange for **shared deployment and administration** of this infrastructure during 24 x 7 operation
- Joint project between LCG sites, experiments and s/w projects
 - Time frame for first service: **deployment in autumn this year**
 - Evolve together with computing models and requirements towards LHC startup
- Held Distributed Database workshop at CERN last December
 - <http://agenda.cern.ch/fullAgenda.php?ida=a044341>

LCG 3D Non-Goals



- 👎 Store all database data
 - Experiments are free to deploy databases and distribute data under their responsibility.
- 👎 Setup a single monolithic distributed database system
 - Given constraints like WAN connections one can not assume that a single synchronously updated database could work and provide sufficient availability.
- 👎 Setup a single vendor system
 - Technology independence and a multi-vendor implementation will be required to minimize the long term risks and to adapt to the different requirements/constraints on different tiers.
- 👎 Impose a CERN centric infrastructure to participating sites
 - CERN is one equal partner of other LCG sites on each tier
- 👎 Decide on an architecture, implementation, new services, policies
 - Produce a technical proposal for all of those to LCG PEB/GDB



WP1 -Data Inventory and Distribution Requirements

- Members are s/w developers from experiments and grid services that use RDBMS data
- Gather data properties (volume, ownership) and requirements
- Integrate access to 3D services into their software

WP2 - Database Service Definition and Implementation

- Members are technology and deployment experts from LCG sites
- Propose a deployment implementation and common deployment procedures

WP3 - Evaluation Tasks

- Short, well defined technology evaluations against the requirements delivered by WP1
- Evaluation are proposed by WP2 (evaluation plan) and typically executed by the people proposing a technology for the service implementation and result in a short evaluation report

Situation on the Application Side



- Databases are used by a growing number of applications in the physics production chain
- Most of these applications are today run centralized but expect to move to distributed deployment
 - for scalability and availability reasons
- This move can be simplified by a common distribution infrastructure
 - but **will not happen by magic**
 - Developers should plan and design new application for distributed deployment
 - Existing (centralised) applications may need to be adapted
- Need to continue to make key applications vendor neutral
 - DB abstraction layers exist or are being implemented in many foundation libraries
 - OGSA-DAI, ODBC, JDBC, ROOT, POOL, ... are steps in this direction
 - Degree of the abstraction achieved varies
 - Still many applications which are only available for one vendor
 - Or have significant schema differences which forbid DB<->DB replications

Site Contacts



- **Established contact to several Tier 1 and Tier 2 sites**
 - Tier 1: ASCC, BNL, CERN, CNAF, FNAL, GridKa, IN2P3, RAL
 - Tier 2: ANL, U Chicago
- **Visited FNAL and BNL around HEPiX**
 - Very useful discussions with experiment developers and database service teams there
- **Regular meetings have started**
 - (Almost) weekly phone meeting back-to-back for
 - Requirement WG
 - Service definition WG
 - Current time (Thursday 16-18)
- **All above sites have expressed interest in project participation**
 - BNL has started to setup Oracle setup
 - Most sites have allocated and installed h/w for participation in the 3D test bed
 - U Chicago agreed to act as Tier 2 in the testbed

Service Definition and Implementation



- **DB Service Discovery**
 - How does a job find a close by replica of the database it needs?
 - Need transparent (re)location of services - eg via a database replica catalog
- **Connectivity, firewalls and connection constraints**
- **Access Control - authentication and authorization**
 - Integration between DB vendor and LCG security models
- **Installation and Configuration**
 - Database server and client installation kits
 - Which database client bindings are required (C, C++, Java(JDBC), Perl, ..) ?
 - Server and client version upgrades (eg security patches)
 - Are transparent upgrades required for critical services?
 - Server administration procedures and tools
 - Need basic agreements to simplify shared administration
 - Monitoring and statistics gathering
- **Backup and Recovery**
 - Backup policy templates, responsible site(s) for a particular data type
 - Acceptable latency for recovery
- **Bottom line: service effort should not be underestimated!**
 - We are rather close to LHC startup and can only afford to propose models that have a good chance of working!
 - Do not just hope for good luck; These services will be a critical part of the experiments' infrastructure and should be handled accordingly!

Database Service Policies



- Several sites have deployment policies in place
 - E.g. FNAL:
 - Staged service levels
 - Development -> Integration -> Production systems
 - Well defined move of new code / DB schema during development process
 - Apps developers and DB experts review and optimize schema before production deployment
- Similar policy proposal prepared for CERN physics database services
 - To avoid recent interference between key production applications of different experiments on shared resources
 - Caused by missing indices, inefficient queries, inadequate hardware resources
 - Storage volume alone is not a sufficient metric to define a database service
 - Need a **reference workload** for each key application to define and optimize the service
 - How many requests from how many clients on how much data are required?
- Especially for distributed DB service this will be essential to avoid surprises on either side

Tier 1 Service Split



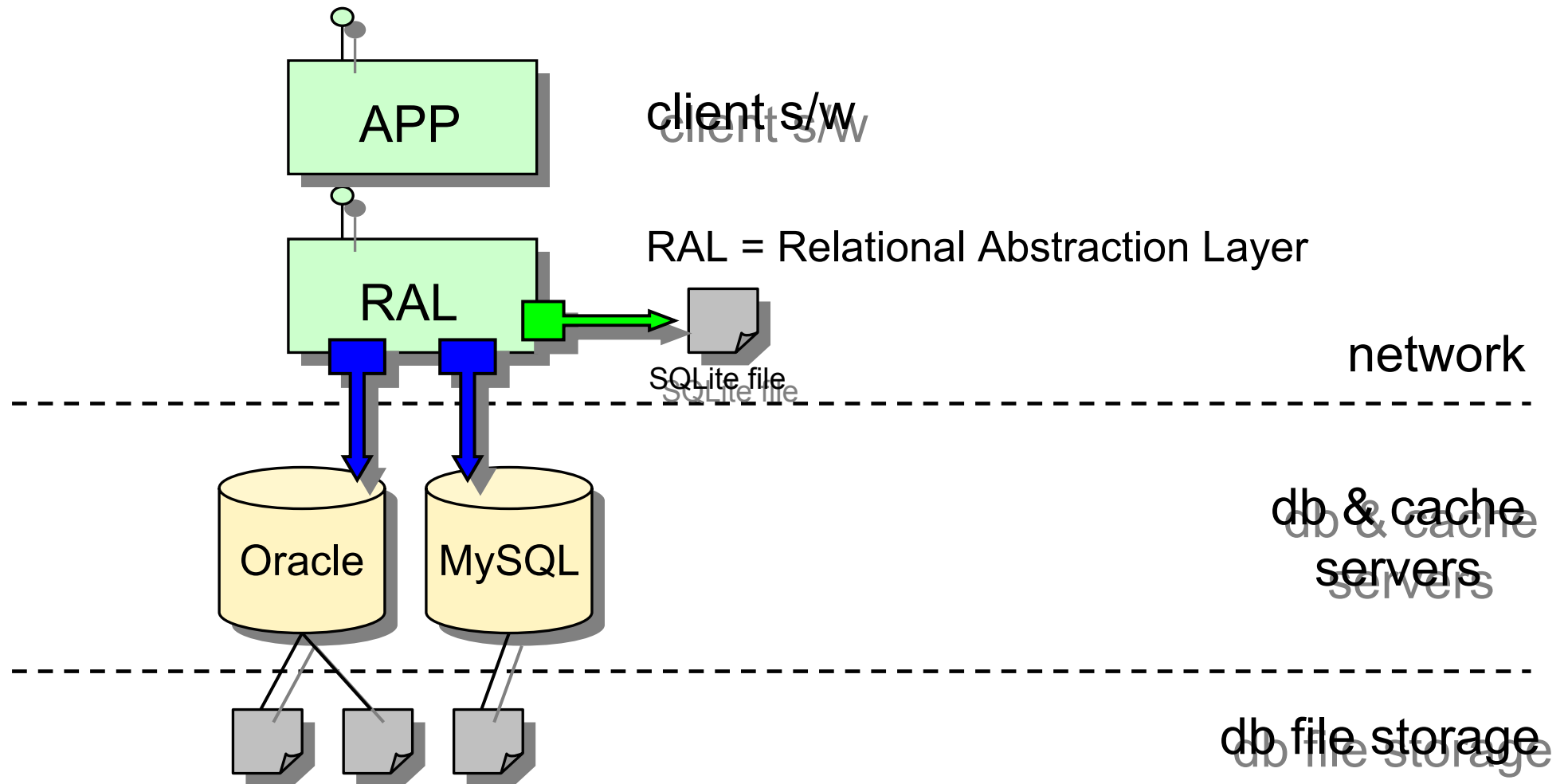
- Discussion only just starting
 - Draft proposal based on input received from FNAL (Anil Kumar)
- Local Services
 - Server installation, bug & security fixes
 - OS patches/upgrades
 - Backup/recovery support
 - Data migration (between db servers)
- Shared Services
 - Db and OS accounts & privileges
 - Storage support (adding more space)
 - Monitoring
 - DB alerts, “killer queries cron job output
 - Host system load, space thresholds
 - Performance problems & optimization
- Site Communication
 - Proposal to setup a shared (web based) Log-Book, mailing lists
 - Need to establish regular DBA meeting
 - eg as part of weekly/bi-weekly 3D meetings

Database Services at LCG Sites Today



- Several tier 1 sites provide Oracle production services for HEP and non-HEP applications
 - Significant deployment experience and well established service exists...
 - ... but can not be changed easily without affecting other site activities
- Tier 2 sites can only provide very limited manpower for database service
 - Part time administration by the same people responsible for fabric
 - Only a simple, constrained database service should be assumed
- MySQL is very popular in the developer community
 - Several applications are bound to MySQL
 - Used for experiment production, though not at very large scale
 - Expected to be deployable with limited db administration resources
- Expect both database flavors to play a role implementing different parts of the LCG infrastructure

Application s/w stack and Distribution Options



Different Distribution Options and their Impact on Application Deployment



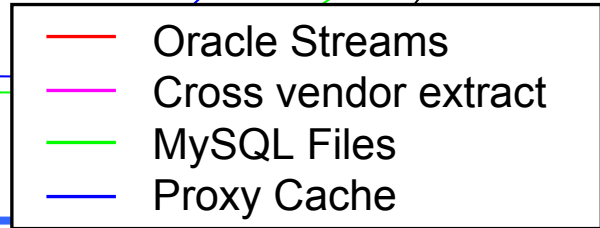
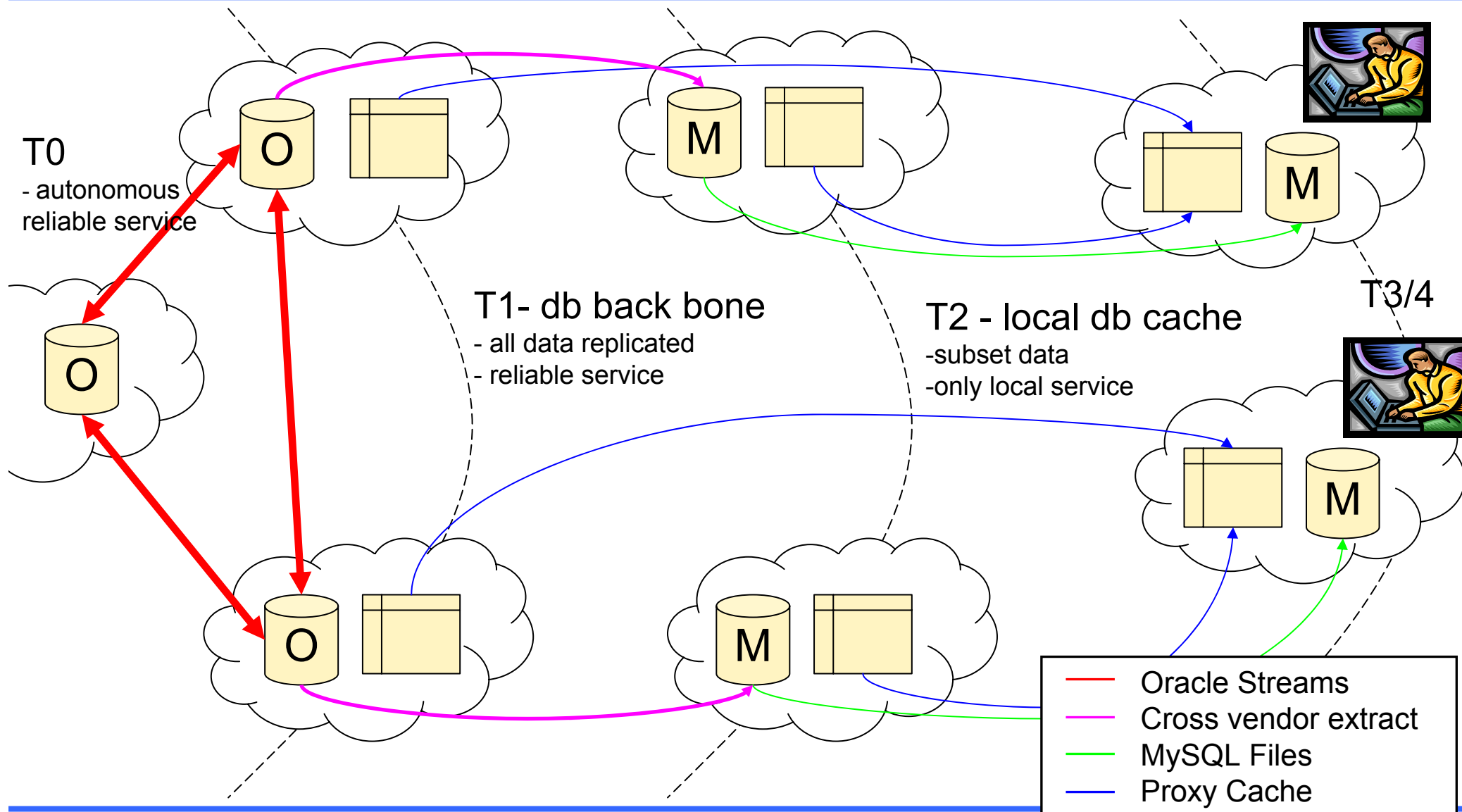
- Database Native **Replication** (eg Oracle Streams)
 - Replication = generic, automated and transactionally consistent data exchange
 - rather than distribution ad hoc, specialized, potentially inconsistent copy
 - Constrains databases vendors
 - Some R&D work starting on cross vendor replication
- Relational Abstraction based **Distribution** (eg ATLAS Octopus)
 - Requires that applications are based on an agreed mapping between different back-ends
 - Needs to be enforced by the abstraction layer or by the application programmer
 - Can cross database vendor boundary
 - One generic process..
 - May run application depending merging/consistency issues
- Application Level **Distribution** (eg POOL File catalogs, old ConditionsDB)
 - Eg using common API or data exchange format shared between different implementations of a component
 - Deployment problems if too many specialized distribution programs need to be maintained and deployed in the distributed environment

Local Database vs. Local Cache



- FNAL developed **FroNtier** system
 - a combination of **http based database access with proxy caches** on the way to the client
 - Performance gains
 - reduced real database access for largely read-only data
 - reduced transfer overhead compared to low level SOAP RPC based approaches
 - Deployment gains
 - Web caches (eg squid) are much simpler to deploy than databases and could remove the need for a local database deployment on some tiers
 - No vendor specific database libraries on the client side
 - “Firewall friendly” tunneling of requests through a single port
- **Expect cache technology to play a significant role**
 - towards the higher tiers which may not have the resources to run a reliable database service

Proposed 3D Service Architecture



Service Integration with Application Software

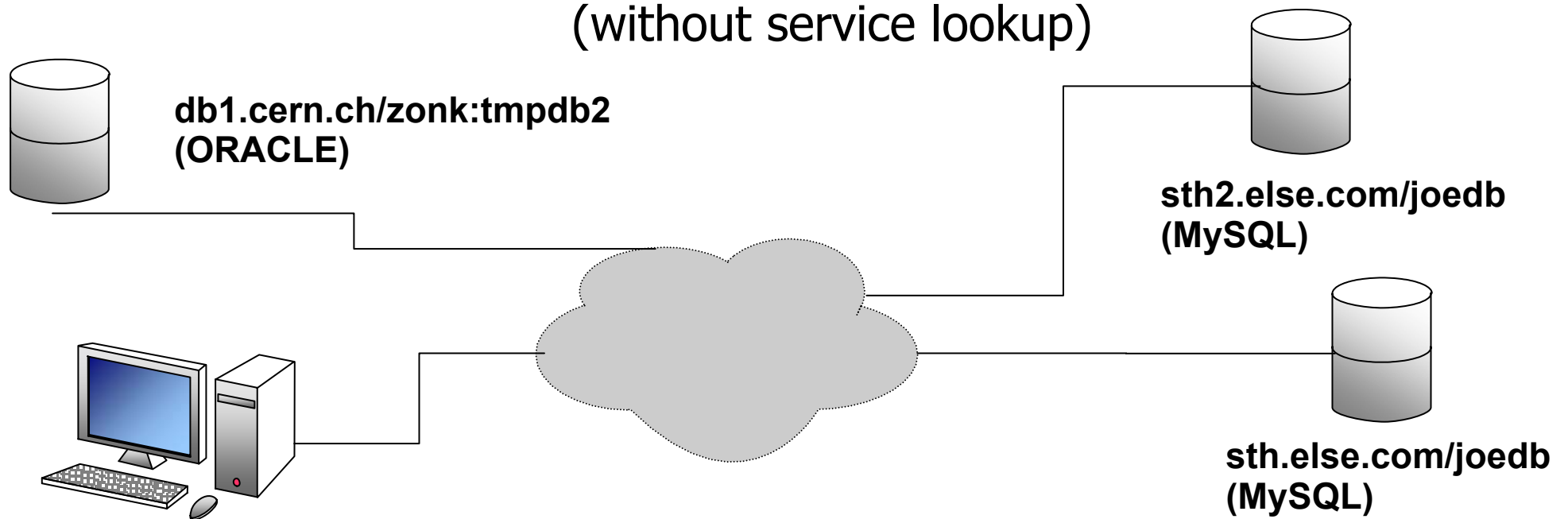


- Distributed database service to be coupled to the user application
 - Plan to use POOL RAL as 3D reference implementation
- How does an application find a suitable database?
 - DB vendor may vary depending on site where job is scheduled
- Database Service Catalog
 - Avoid embedding physical connection strings in code
 - ..or spreading them as part of configuration file copies
 - Allow a local service to be transparently relocated to another machine
- 3D prototype based on POOL file catalog
 - supports network disconnected lookup based on XML files
 - Final implementation may still change
 - Eg if a suitable grid-service has been identified

Database Service Lookup



Usual case (without service lookup)



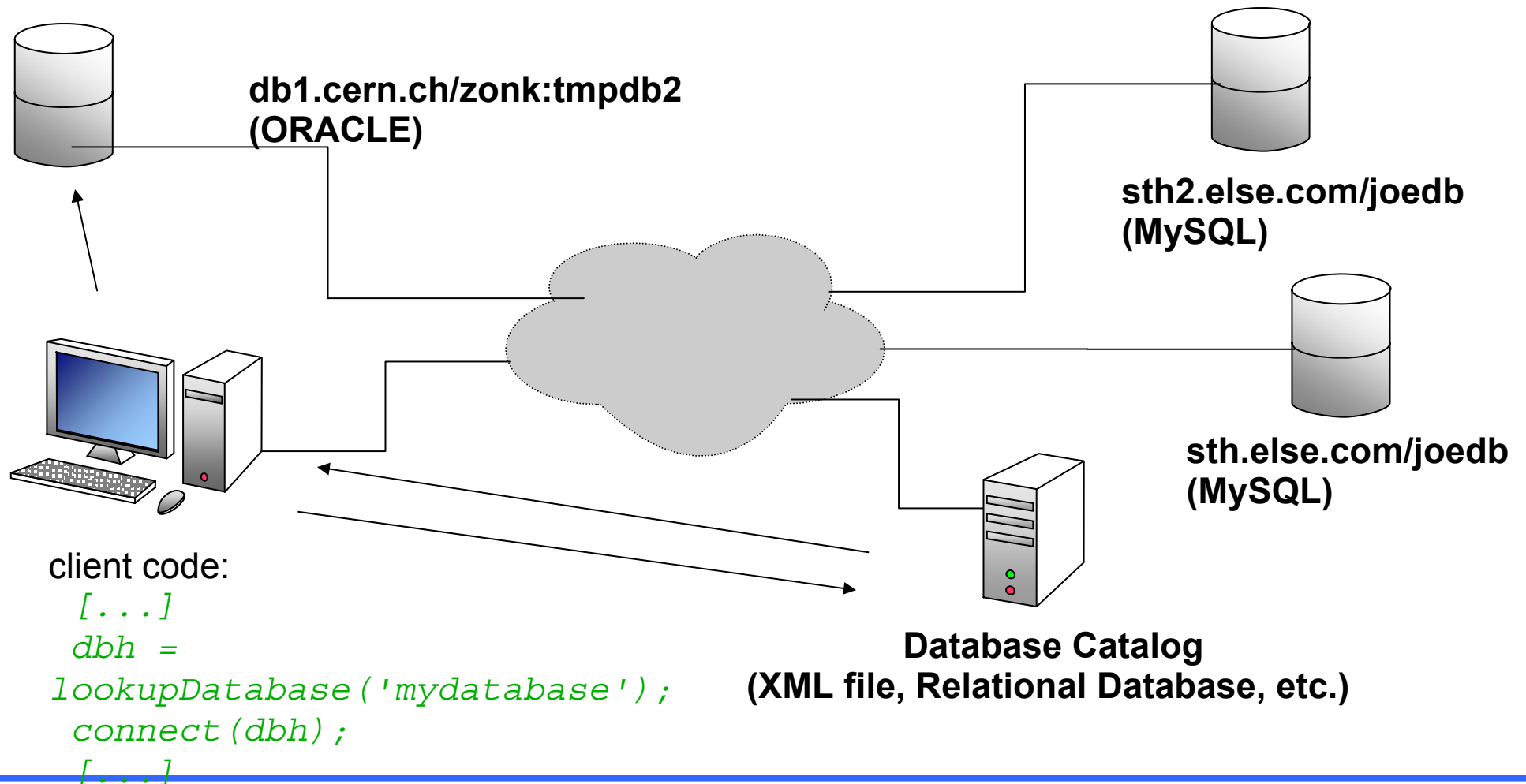
client code:

```
[...]  
connect ('mysql://sth2.else.com/joedb');  
[...]
```

Database Service Lookup



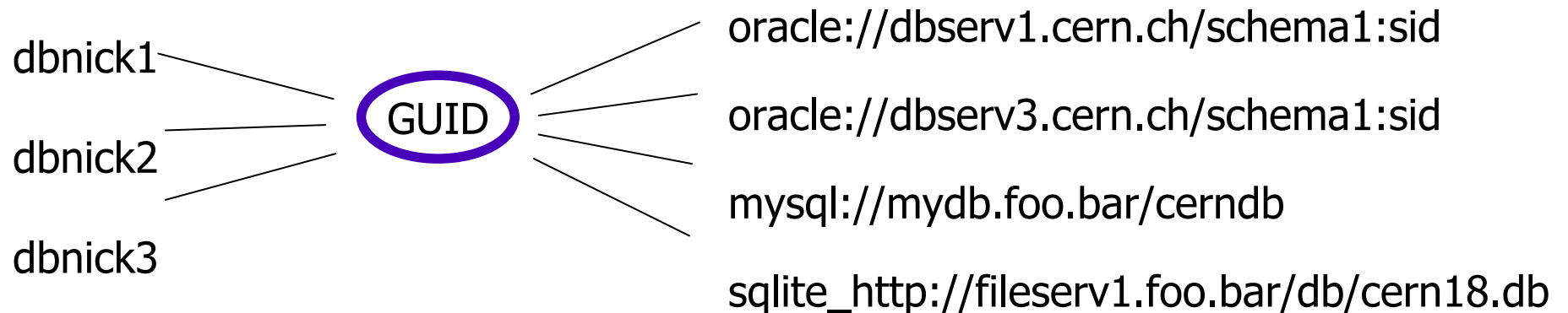
With Service Lookup



Database Catalog



- Terms
 - Logical Database Name
 - Physical Database (connection string)
 - Internally: GUID (Globally Unique Identifier)
- Prototype developed reusing POOL file catalog components



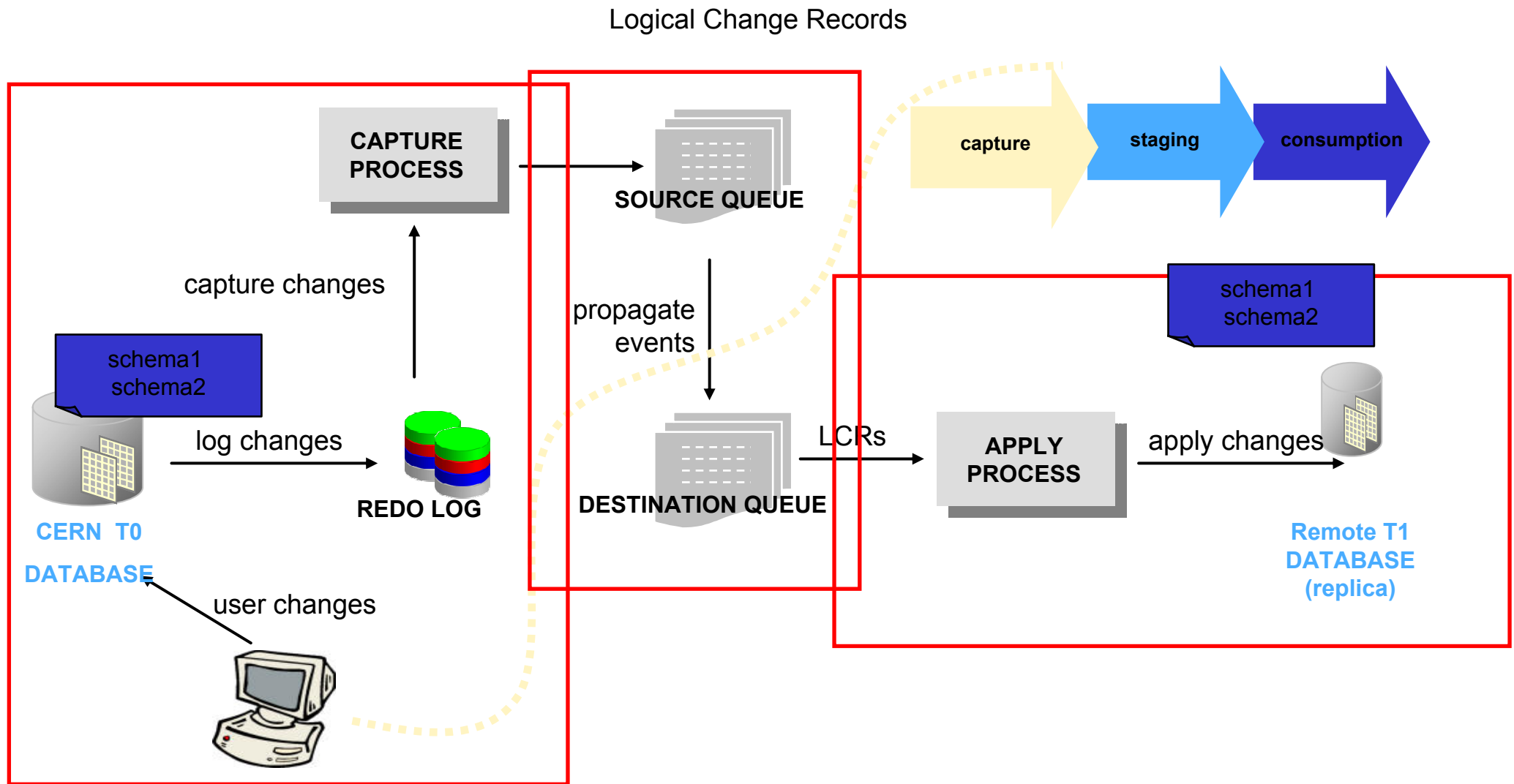
▶ **No usernames or passwords stored!**

More Application Integration..

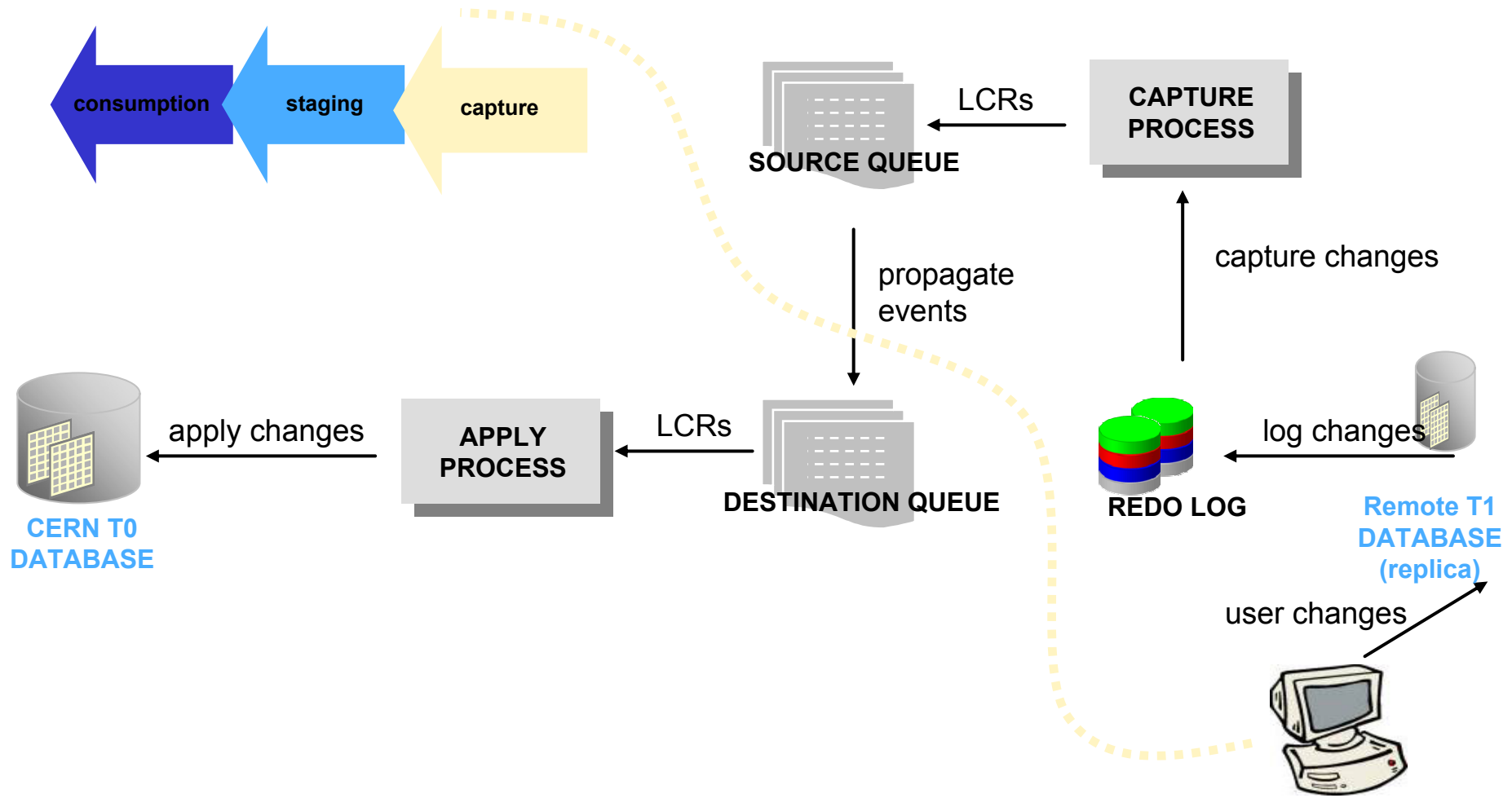


- **Connection Pooling and Error Handling**
 - ATLAS developments being integrated into RAL
- **Authentication and Authorization**
 - Mapping between grid certificate and database role
 - Need to support roles in experiments computing model
 - Plan to evaluate/integrate with package developed by G.Ganis for ROOT and xrootd
- **Client Diagnostics**
 - Collect timing of top queries of the current application in RAL
 - Evaluate package developed at FNAL to aggregate diagnostics across many clients
- **Most of the above are useful for local or centralized database applications too**

STREAMS Architecture (I)



STREAMS Architecture (I)



LCG 3D Testbed



- Oracle 10g server
 - Install kits and documentation are provided for test bed sites
 - CERN can not offer offsite support though
 - At least 100GB storage
 - Application and reference load packaged by CERN / experiments
- FroNtier installation
 - Just one server plus squid installations at other sites?
 - Need squid package and install instructions (FNAL?)
 - Need a test server at CERN or FNAL
- Worker nodes to run reference load
- Oracle Enterprise Manager installation for test bed administration and diagnostic

3D Data Inventory



- Collect and maintain a catalog of main RDBMS data types
- Experiments and grid s/w providers fill a table for each data type which is candidate for storage and replication via the 3D service
 - Basic storage properties
 - Data description, expected volume on T0/1/2 in 2005 (and evolution)
 - Ownership model: read-only, single user update, single site update, concurrent update
 - Replication/Caching properties
 - Replication model: site local, all t1, sliced t1, all t2, sliced t2 ...
 - Consistency/Latency: how quickly do changes need to reach other sites/tiers
 - Application constraints: DB vendor and DB version constraints
 - Reliability and Availability requirements
 - Essential for whole grid operation, for site operation, for experiment production,
 - Backup and Recovery policy
 - acceptable time to recover, location of backup(s)

Requirement Gathering



- All participating experiments have signed up their representatives
 - Rather 3 than 2 names
 - Tough job, as survey inside experiment crosses many boundaries!
- Started with a simple spreadsheet capturing the various applications
 - Grouped by application
 - One line per replication and Tier (distribution step)
 - Two main patterns
 - Experiment data: data fan-out - $T_0 \rightarrow T_n$
 - Grid services: data consolidation - $T_n \rightarrow T_0$
 - But some exceptions which need to be documented...
- Aim to collect complete set of requirements for database services
 - Eg also online data or data which is stored locally but never leaves a tier
 - Needed to properly size the h/w at each site

Preliminary Summary for 2005



- Applications mentioned so far
 - FileCatalog, Conditions, Geometry, Bookkeeping, Physics Meta Data, Collections, Grid Monitoring, TransferDB
 - Suspect that several smaller(?) applications are still missing
- Total volume per experiment: 50-500 GB
 - Error bars likely as big as the current spread
 - Significant flexibility/headroom required for service side!
- Number of applications to be supported: O(5)
 - Some still not existing or bound to MySQL at the moment
 - Distributed applications use either RAL or ODBC based implementation
- Distributed Data becomes read-only down from T0
 - Conservative approach for first service deployment
 - Current model **does not require multi-master replication**
- **Please make sure your experiment representatives know about your application requirements!**

Summary



- A distributed database infrastructure promises to provide scalability and availability for database applications at LHC
 - The LCG 3D project has been started as joint project between experiments and LCG sites to coordinate the definition of this service
 - Several distribution options are available for planning/design of new applications
- **DB Service Definition**
 - Very relevant deployment/development experience from RUN2 @ FNAL
 - Service task split and application validation policy proposals are firming up
 - Oracle 10g based replication test-bed expands to first T1
- **Several new DB applications will be developed this year**
 - Expect high support load on database services (especially at Tier 0) for 2005
 - Starting to think early about concrete deployment and distribution models is key for successful production
- **Please contact us to discuss the requirements and distribution options for your application with you.**