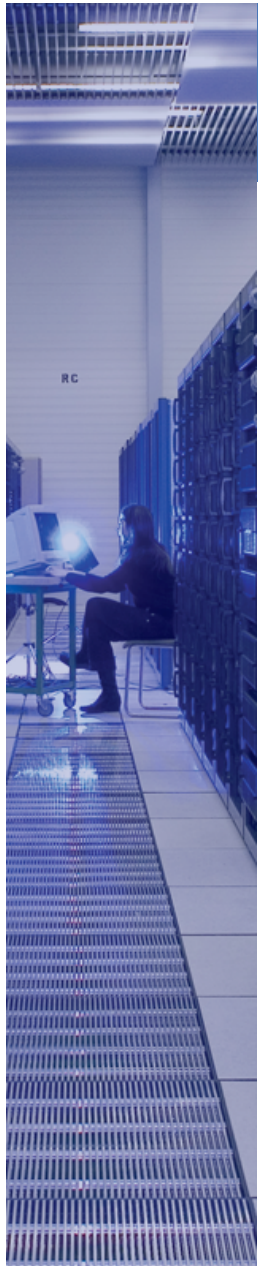


# Worldwide distribution of experimental physics data using Oracle Streams

Eva Dafonte Pérez

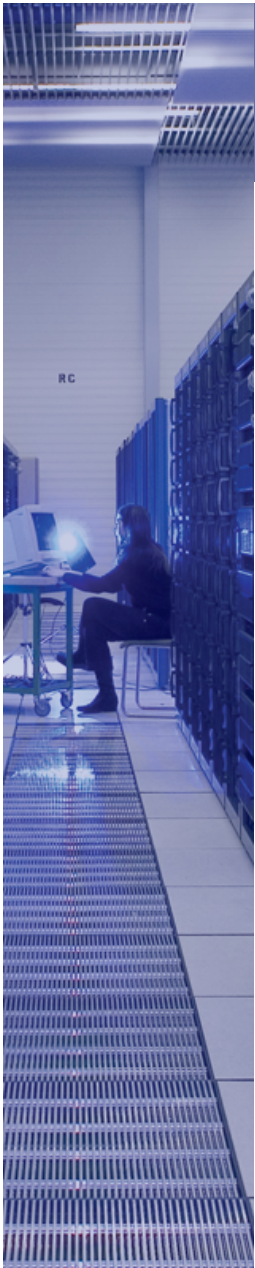
Database Administrator @CERN



# Outline

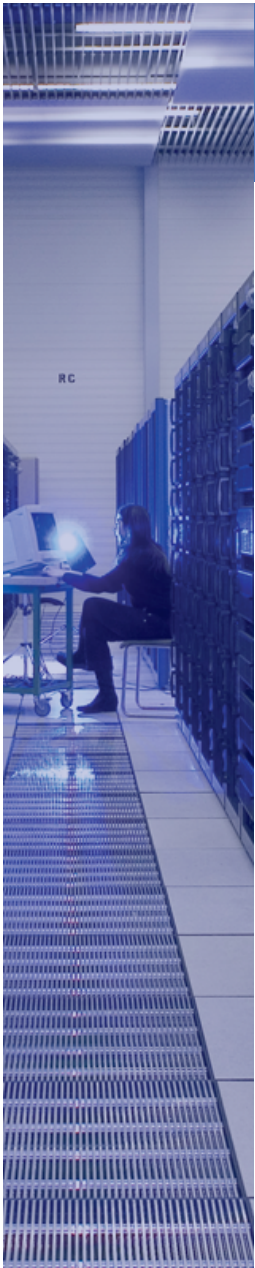
- CERN and LHC Overview
- Oracle Streams Replication
- Replication Performance
- Optimizations: Downstream Capture, Split and Merge, Network, Rules and Flow Control
- Periodic Maintenance
- Lessons Learned
- Tips and Tricks
- Streams Bugs and Patches
- Scalable Resynchronization
- 3D Streams Monitor
- New 11g Streams Features
- Streams Setups Examples
- Summary



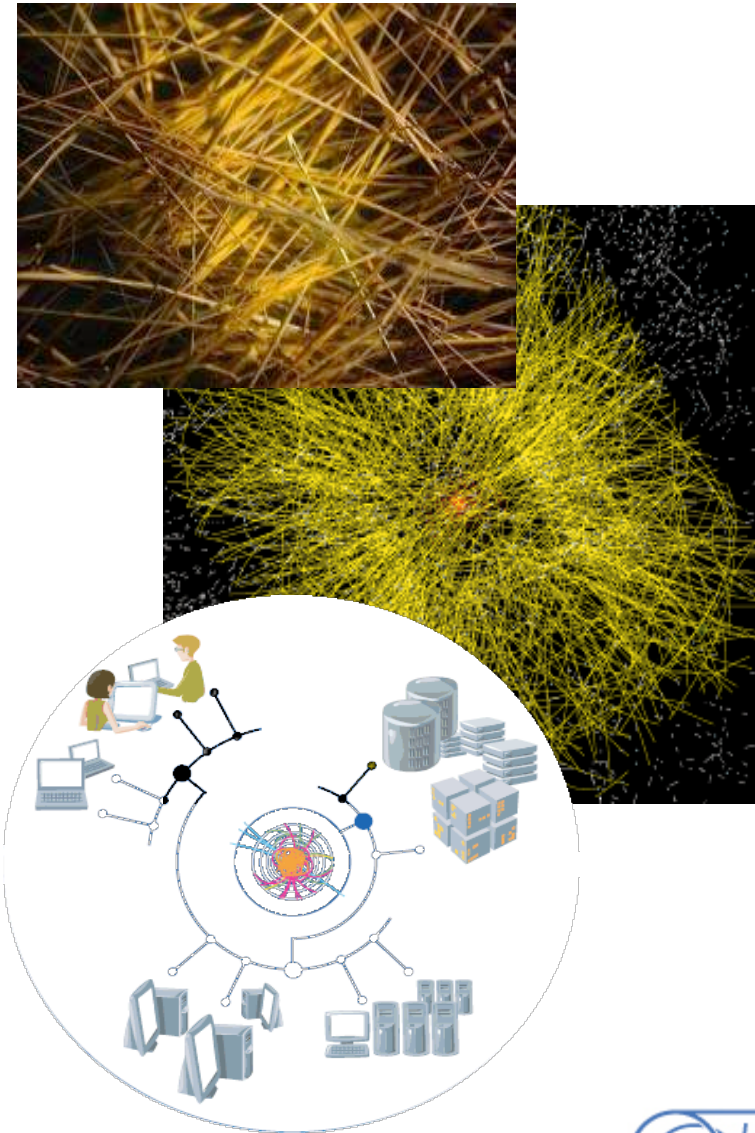


- European Organization for Nuclear Research
  - world's largest centre for scientific research
  - founded in 1954
  - mission: finding out what the Universe is made of and how it works
- LHC, Large Hadron Collider
  - particle accelerator used to study the smallest known particles
  - 27 km ring, spans the border between Switzerland and France about 100 m underground
  - will recreate the conditions just after the Big Bang

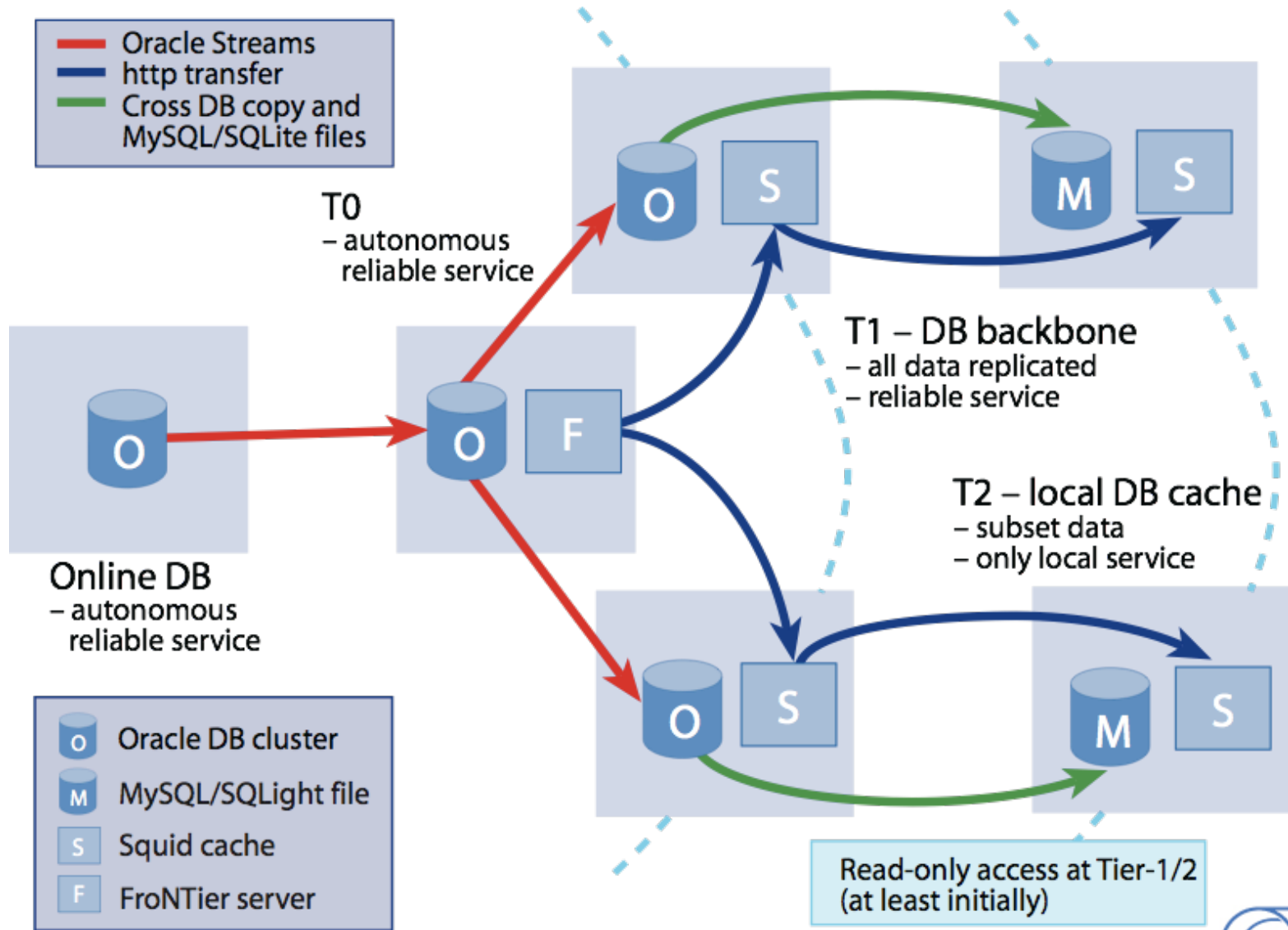
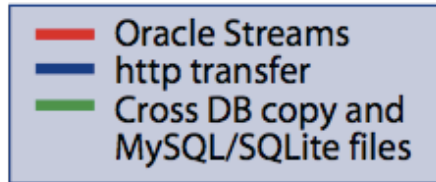
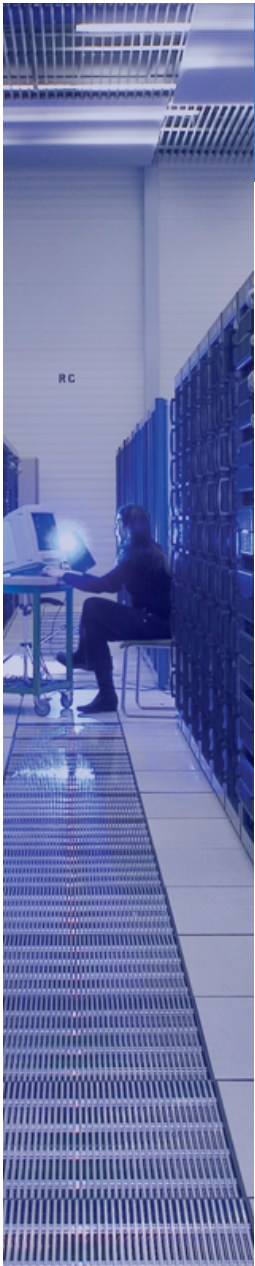
# The LHC Computing Challenge



- Data volume
  - high rate x large number of channels x 4 experiments
  - **15 PetaBytes of new data each year stored**
  - **much more data discarded during multi-level filtering before storage**
- Compute power
  - event complexity x Nb. events x thousands users
  - **100 k of today's fastest CPUs**
- Worldwide analysis & funding
  - computing funding locally in major regions & countries
  - efficient analysis everywhere
  - **GRID technology**

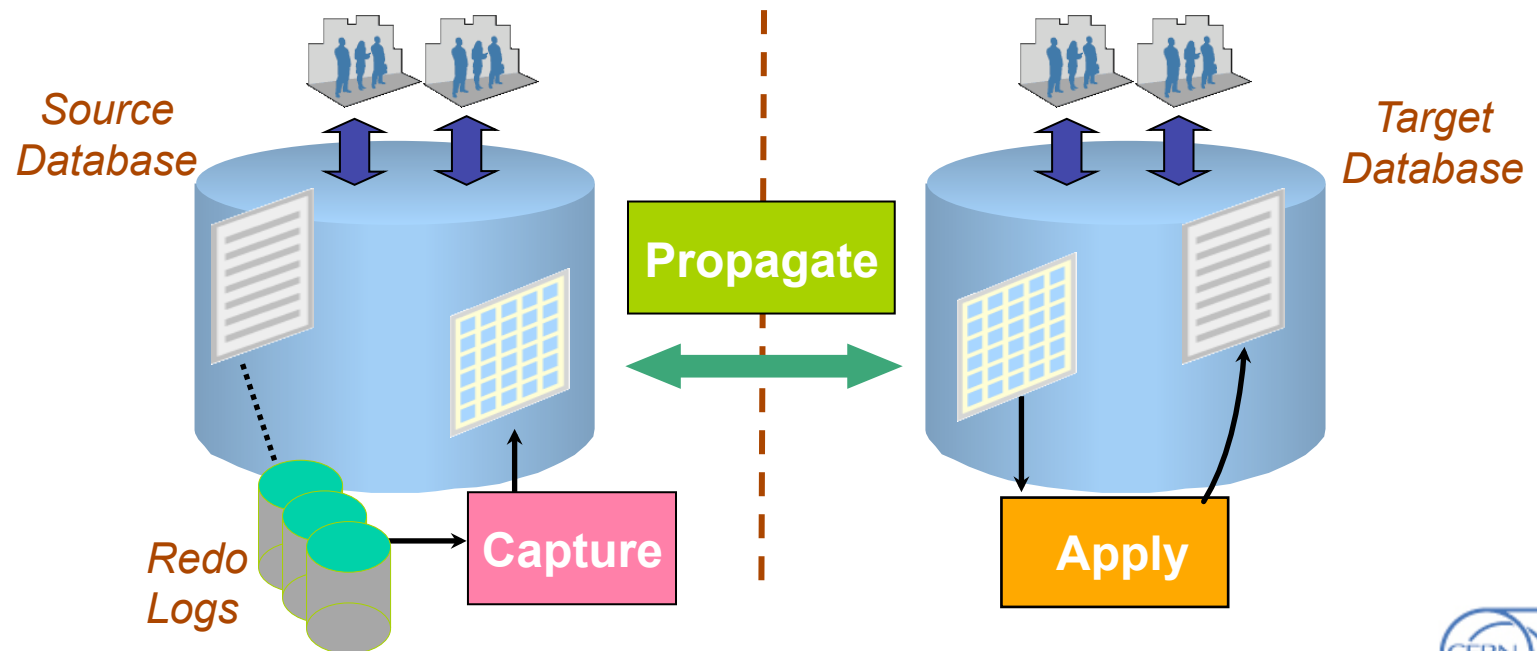


# Distributed Service Architecture



# Oracle Streams Replication

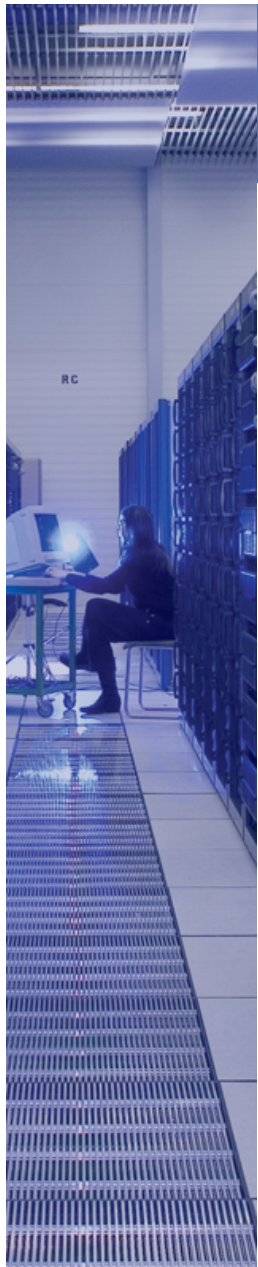
- Technology for **sharing information** between databases
- Database changes captured from the redo-log and propagated asynchronously as Logical Change Records (LCRs)



- The atomic unit is the change record: LCR
- LCRs can vary widely in size
  - Throughput is not a fixed measure
- Capture performance:
  - Read changes from the redo
    - from redo log buffer (memory - much faster)
    - from archive log files (disk)
  - Convert changes into LCRs
    - depends on the LCR size and number of columns
  - Enqueue the LCRs
    - concurrent access to the data structure can be costly



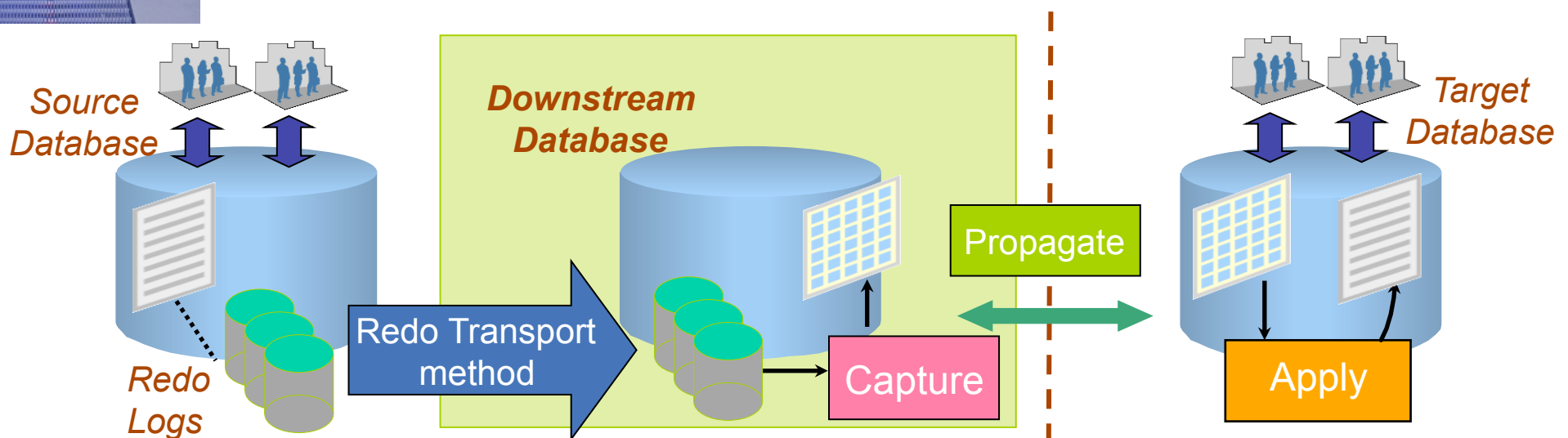
- Propagation performance:
  - Browse LCRs
  - Transmit LCRs over the network
  - Remove LCRs from the queue
    - Done in separate process to avoid any impact
- Apply performance:
  - Browse LCRs
  - Execute LCRs
    - Manipulate the database is slower than the redo generation
    - Execute LCRs serially => apply cannot keep up with the redo generation rate
  - Remove LCRs from the queue



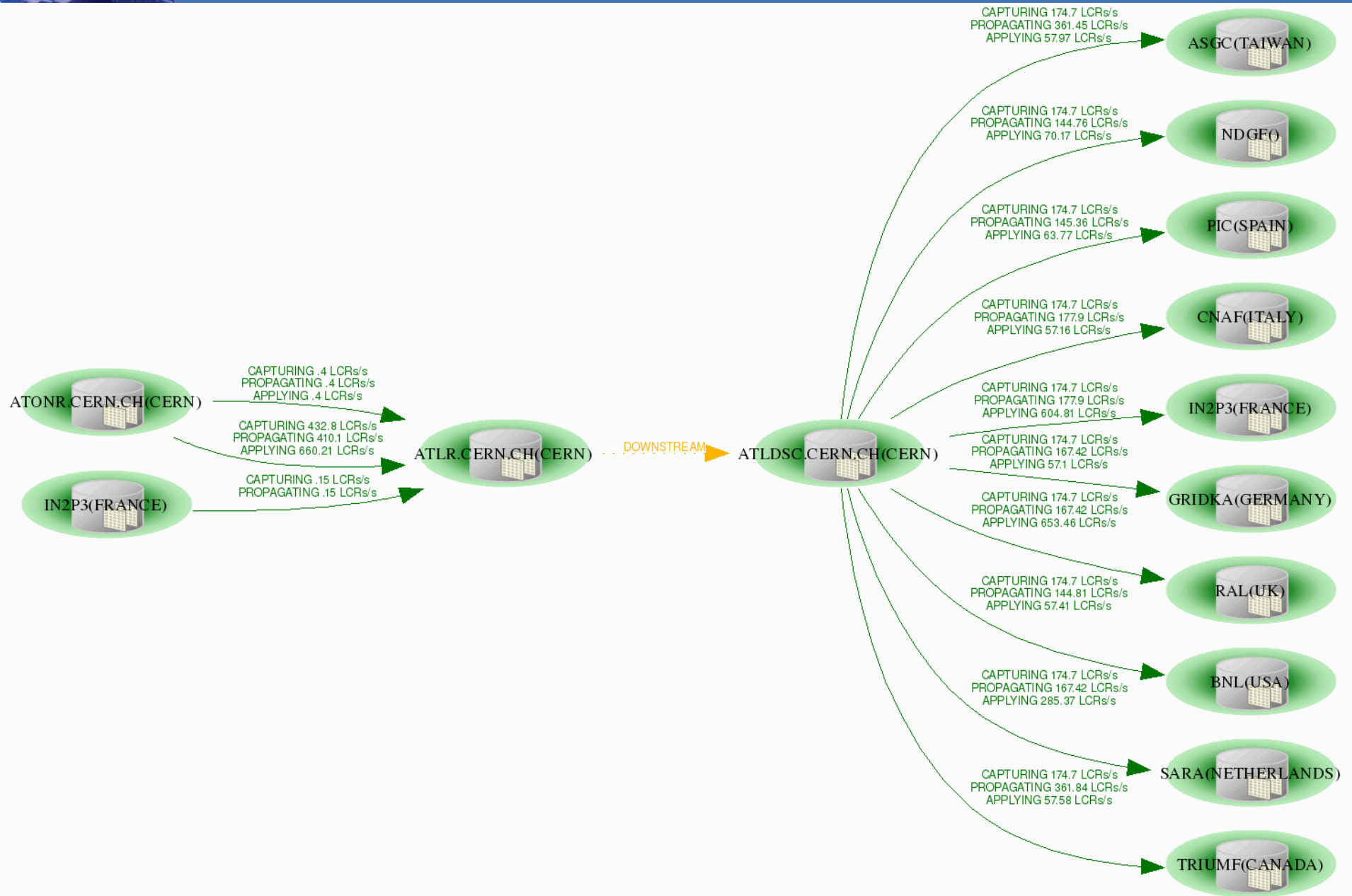


# Downstream Capture

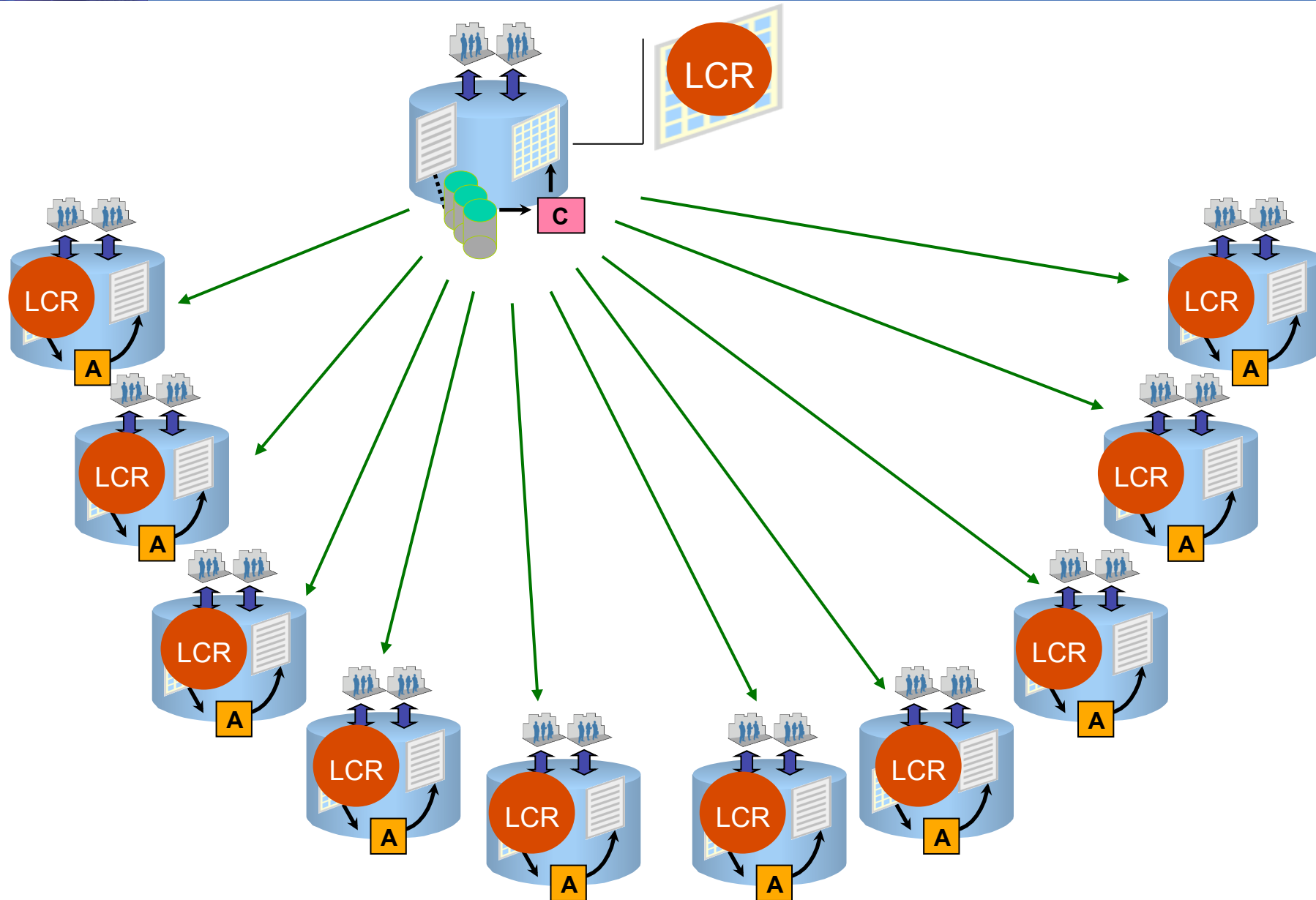
- **Downstream capture** to de-couple Tier 0 production databases from destination or network problems
  - source database availability is highest priority
- Optimizing **redo log retention** on downstream database to allow for sufficient re-synchronisation window
  - we use 5 days retention to avoid tape access
- Dump fresh copy of dictionary to redo periodically
- 10.2 Streams recommendations (metalink note 418755)



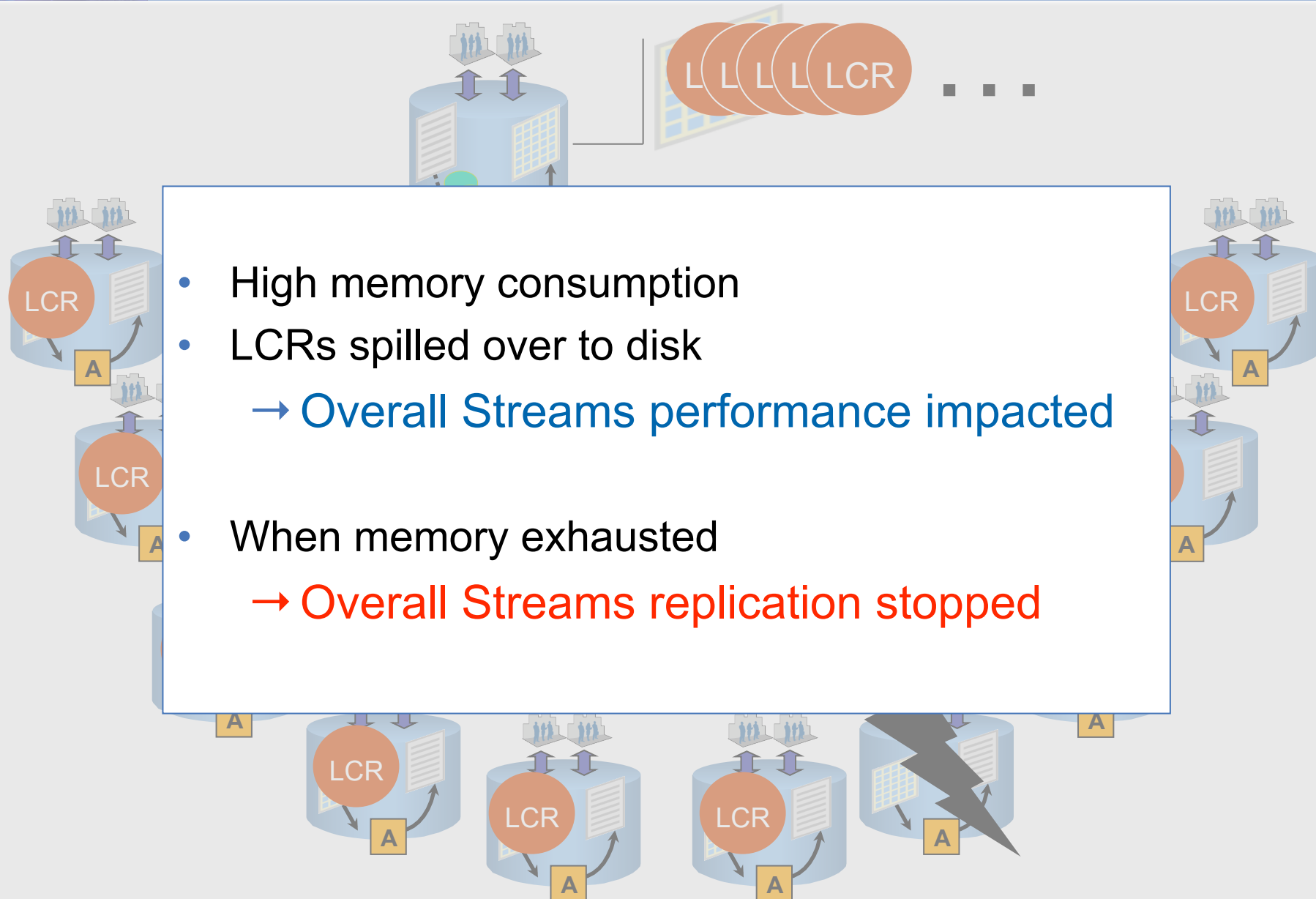
# Streams Setup Example: ATLAS



# Split & Merge: Motivation



# Split & Merge: Motivation



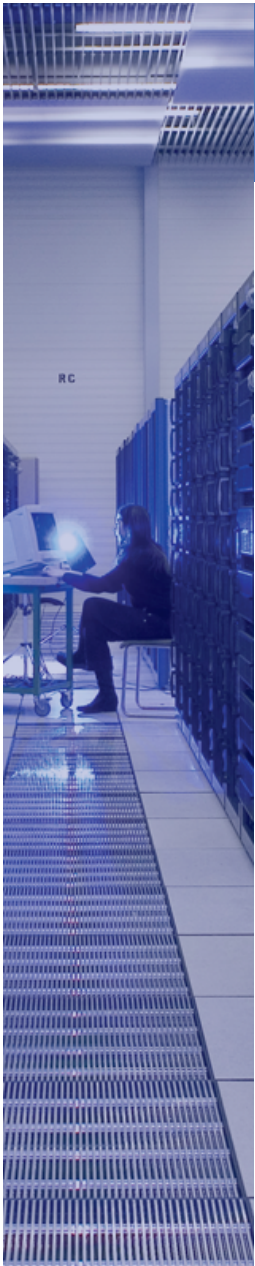
# Split & Merge

*in collaboration with Patricia McElroy  
Principal Product Manager Distributed Systems/Replication - Oracle*

- Objective: **isolate replicas against each other**
  - Split
    - (original) Streams setup for “good” sites
      - drop propagation job/s to “bad” site/s
        - spilled LCRs are removed from the capture queue
    - (new) Streams setup for “bad” site/s
      - new capture queue
      - clone capture process and propagation job/s
    - does not require any change on the destination site/s
  - Merge
    - move back the propagation job/s to the original setup
    - clean up additional Streams processes and queue
    - does not require any change on the destination site/s



# Split & Merge: Details



- **Split:**

```
SQL> exec split ('STRM_PROP_A', 'STRM_CAP_CL', 'STRMQ_CL', 'STRM_PROP_CL');  
exec resynchronize_site ('STRMTEST.CERN.CH', 'STRM_CAP_CL',  
'STRMQ_CL', 1, 2, 'STRM_PROP_CL', 'STRMQ_A_AP', 'RULESET$_18', '');
```

- gather cloning information:

- **capture:**

- rule set name

- start\_scn = last applied message scn @target

- first\_scn = previous dictionary build < start\_scn

- **propagation:**

- rule set name

- target queue name and db link

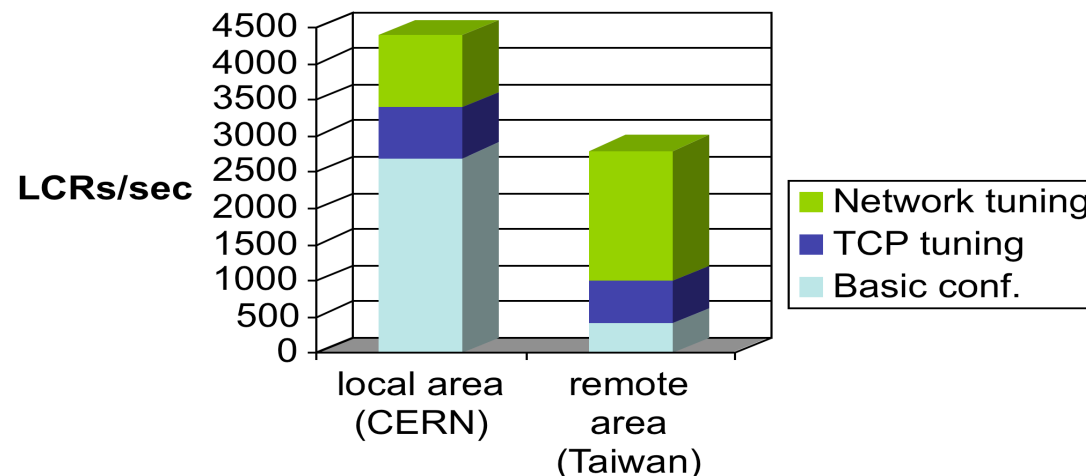
- **Merge:**

```
SQL> exec merge('STRM_CAP_SA', 'STRM_CAP_CL', 'STRM_PROP_A', 'STRM_PROP_CL');
```

- select the minimum required checkpoint scn between the 2 capture processes

- recover original propagation

- TCP and Network tuning
  - adjust system max TCP buffer (/etc/sysctl.conf)
  - parameters to reinforce the TCP tuning
    - DEFAULT\_SDU\_SIZE=32767
    - RECV\_BUF\_SIZE and SEND\_BUF\_SIZE
      - Optimal:  $3 * \text{Bandwidth Delay Product}$
- Reduce the Oracle Streams acknowledgements
  - alter system set events '26749 trace name context forever, level 2';



# Streams Rules

- Used to control which information to share
- Rules on the capture side caused more overhead than on the propagation side
- **Avoid Oracle Streams complex rules**

## Complex Rule

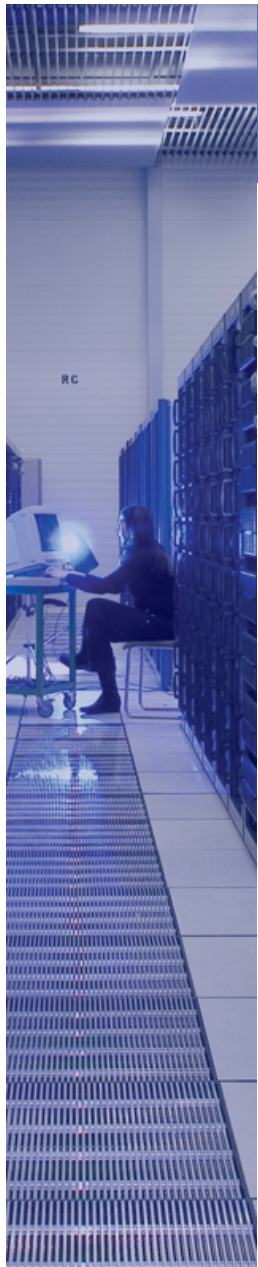
```
condition => '( SUBSTR(:ddl.get_object_name(),1,7) IN ("COMP200", "OFLP200", "CMCP200",
"TMCP200", "TBDP200", "STRM200")
OR SUBSTR (:ddl.get_base_table_name(),1,7) IN ("COMP200", "OFLP200", "CMCP200",
"TMCP200", "TBDP200", "STRM200") )'
```

### Avoid complex rules:

- LIKE
- Functions
- NOT

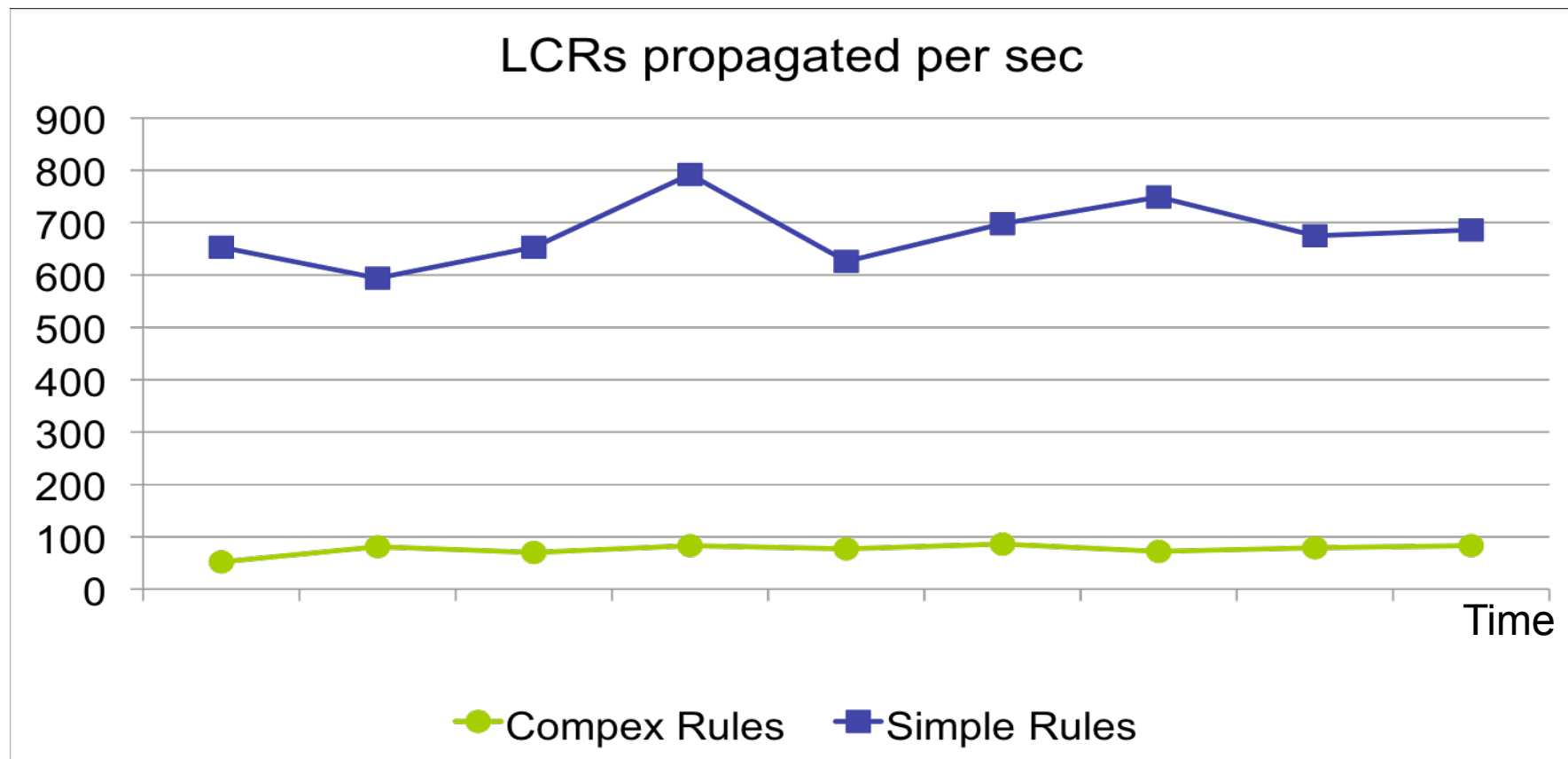
## Simple Rule

```
condition => '(((ddl.get_object_name() >= "STRM200_A" and :ddl.get_object_name() <= "STRM200_Z") OR
(ddl.get_base_table_name() >= "STRM200_A" and :ddl.get_base_table_name() <= "STRM200_Z"))
OR ((ddl.get_object_name() >= "OFLP200_A" and :ddl.get_object_name() <= "OFLP200_Z") OR
(:ddl.get_base_table_name() >= "OFLP200_A" and :ddl.get_base_table_name() <= "OFLP200_Z"))'
```

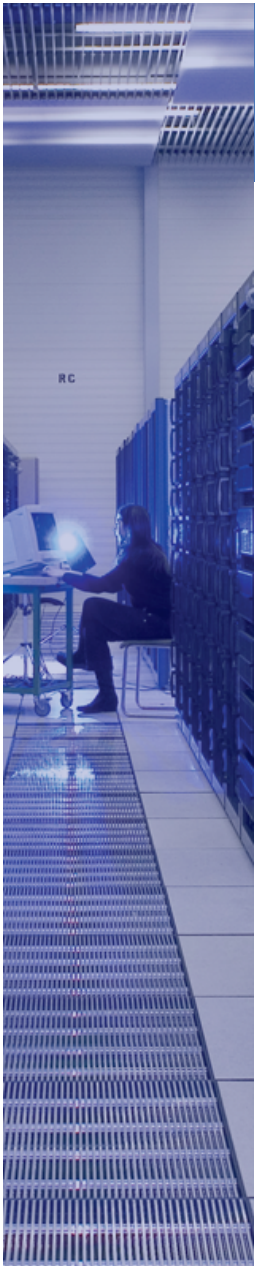




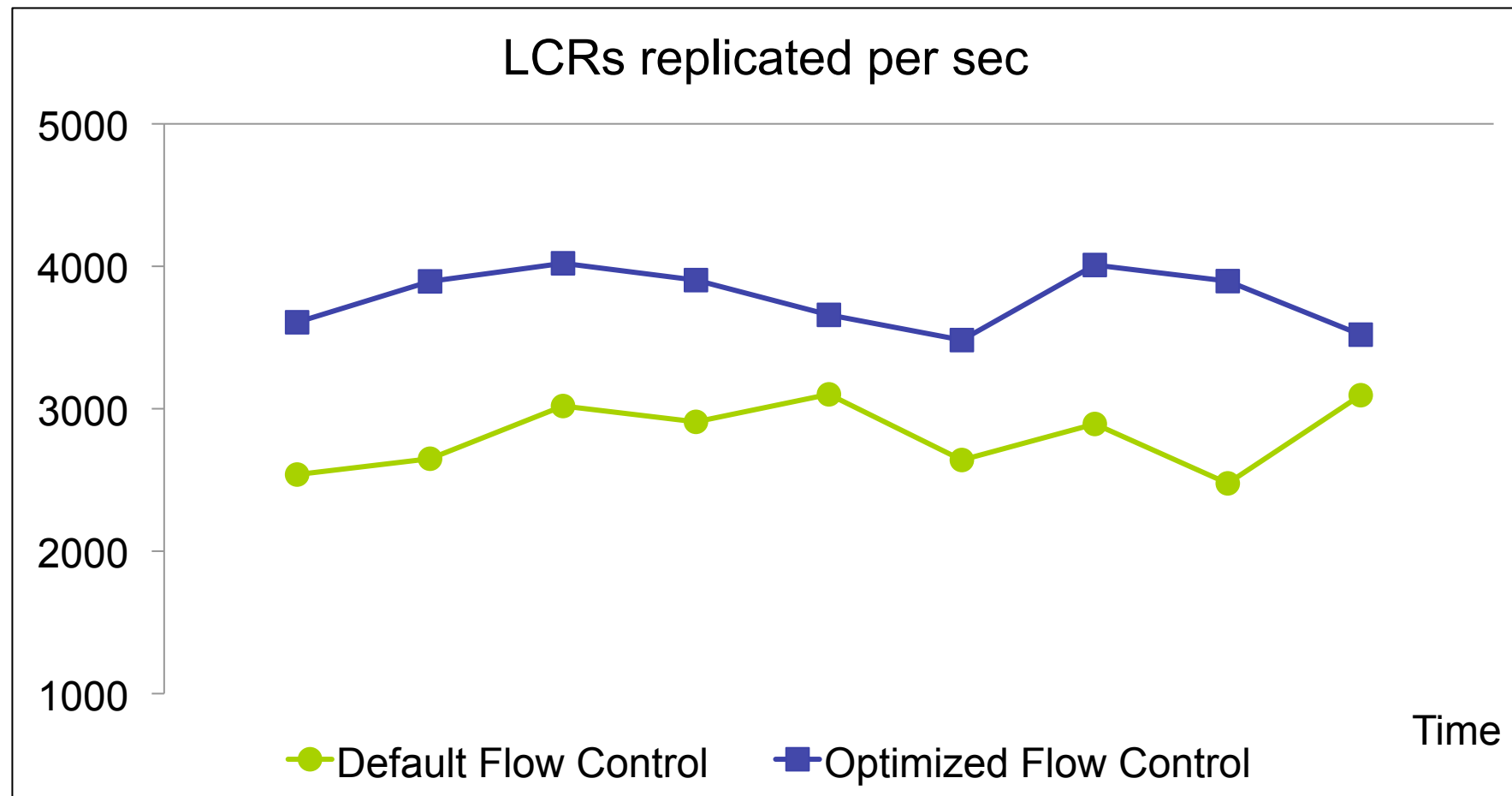
- Example: ATLAS Streams Replication
  - rules defined to filter tables by prefix

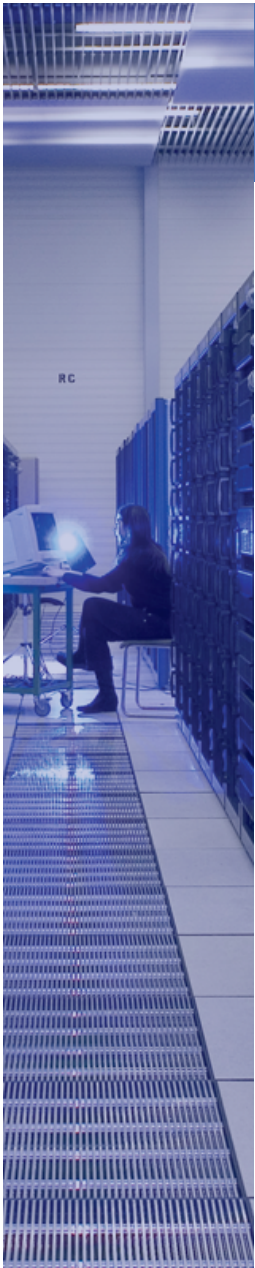


- By default, flow control kicks when the number of messages is larger than the threshold
  - Buffered publisher: 5000
  - Capture publisher: 15000
- **Manipulate default behavior**
- 10.2.0.3 + Patch 5093060 = 2 new events
  - 10867: controls threshold for any buffered message publisher
  - 10868: controls threshold for capture publisher
- 10.2.0.4 = 2 new hidden parameters
  - “\_capture\_publisher\_flow\_control\_threshold”
  - “\_buffered\_publisher\_flow\_control\_threshold”

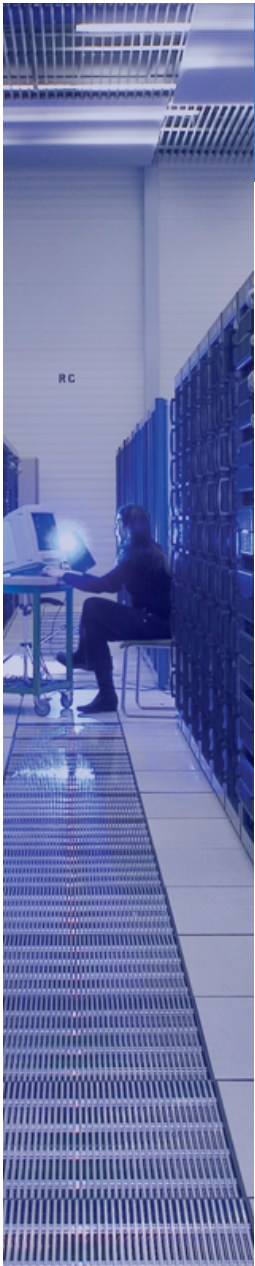


- Example: ATLAS PVSS Streams Replication



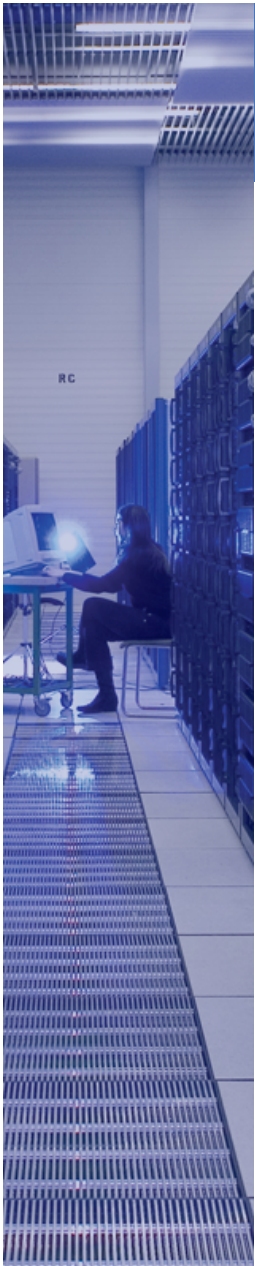


- **Dump fresh copy of Dictionary redo**
  - reduces the amount of logs to be processed in case of additional process creation
- **Reduce high watermark of AQ objects**
  - maintain enqueue/dequeue performance
  - reduce QMON CPU usage
  - metalink note 267137.1
- **Shrink Logminer checkpoint table**
  - improves capture performance
  - metalink note 429599.1
- **Review the list of specific Streams patches**
  - metalink note 437838.1



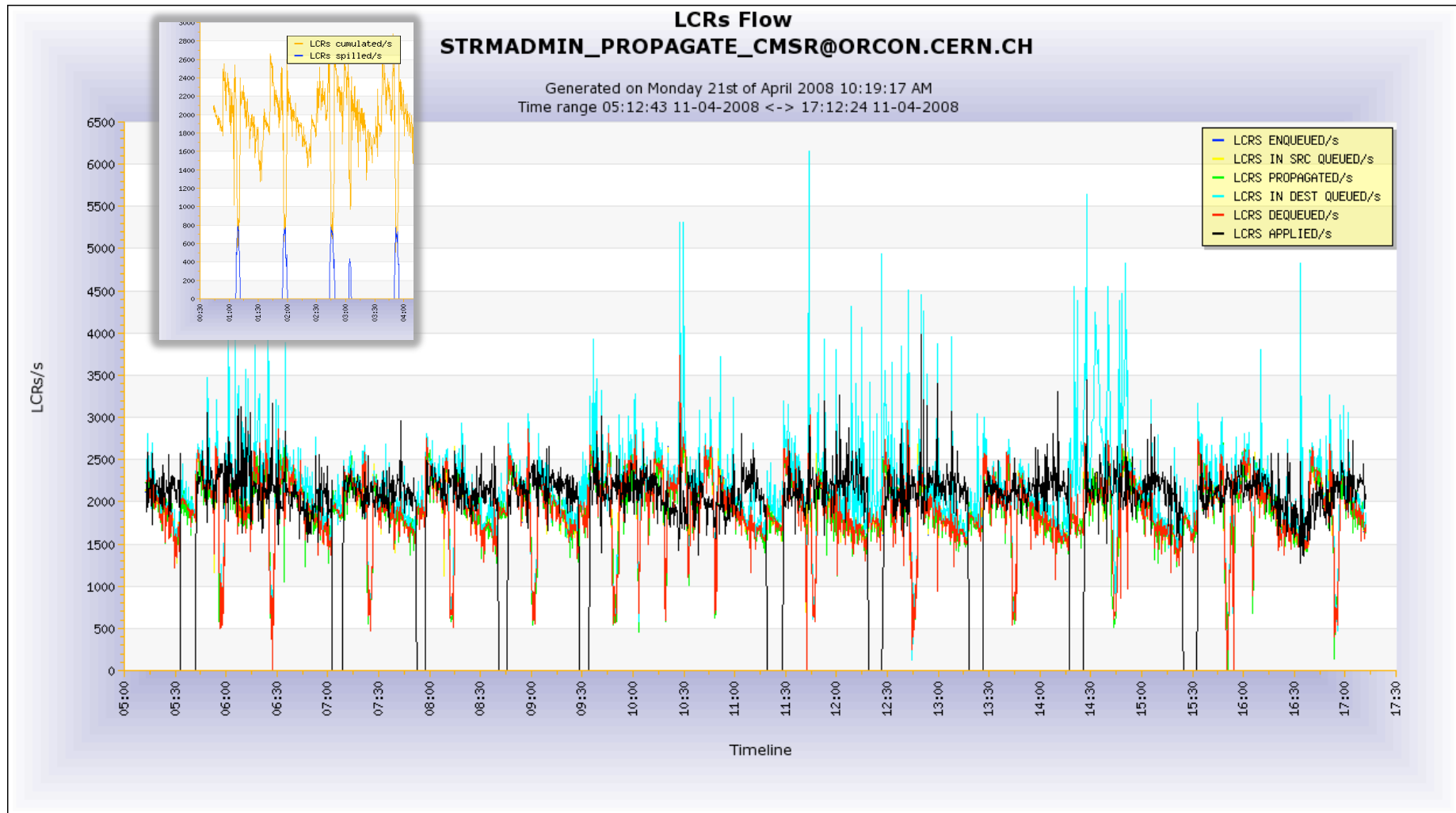
- SQL bulk operations (at the source db)
  - may map to many elementary operations at the destination side
  - need to control source rates to avoid overloading
- Batch processing
  - minimize the performance impact **using Streams tags**
  - avoid changes being captured, then run same batch load on all destination
- System generated names
  - do not allow system generated names for constraints and indexes
  - modifications will fail at the replicated site
  - storage clauses also may cause some issues if the target sites are not identical

- Replication of “grant” operations
  - grants on views, PL/SQL procedures, functions and packages are NOT replicated
  - grantee must exist at all destinations
- Long transactions (non-frequent commits)
  - Total number of outstanding LCRs is too large
  - LCRs are in memory too long
    - LCRs are spilled over to disk
    - Apply performance is impacted
  - All LCRs in a single transaction must be applied by one apply server
    - Parallel servers cannot be used efficiently
  - Too many unbrowsed messages enables flow control
    - Streams processes are paused



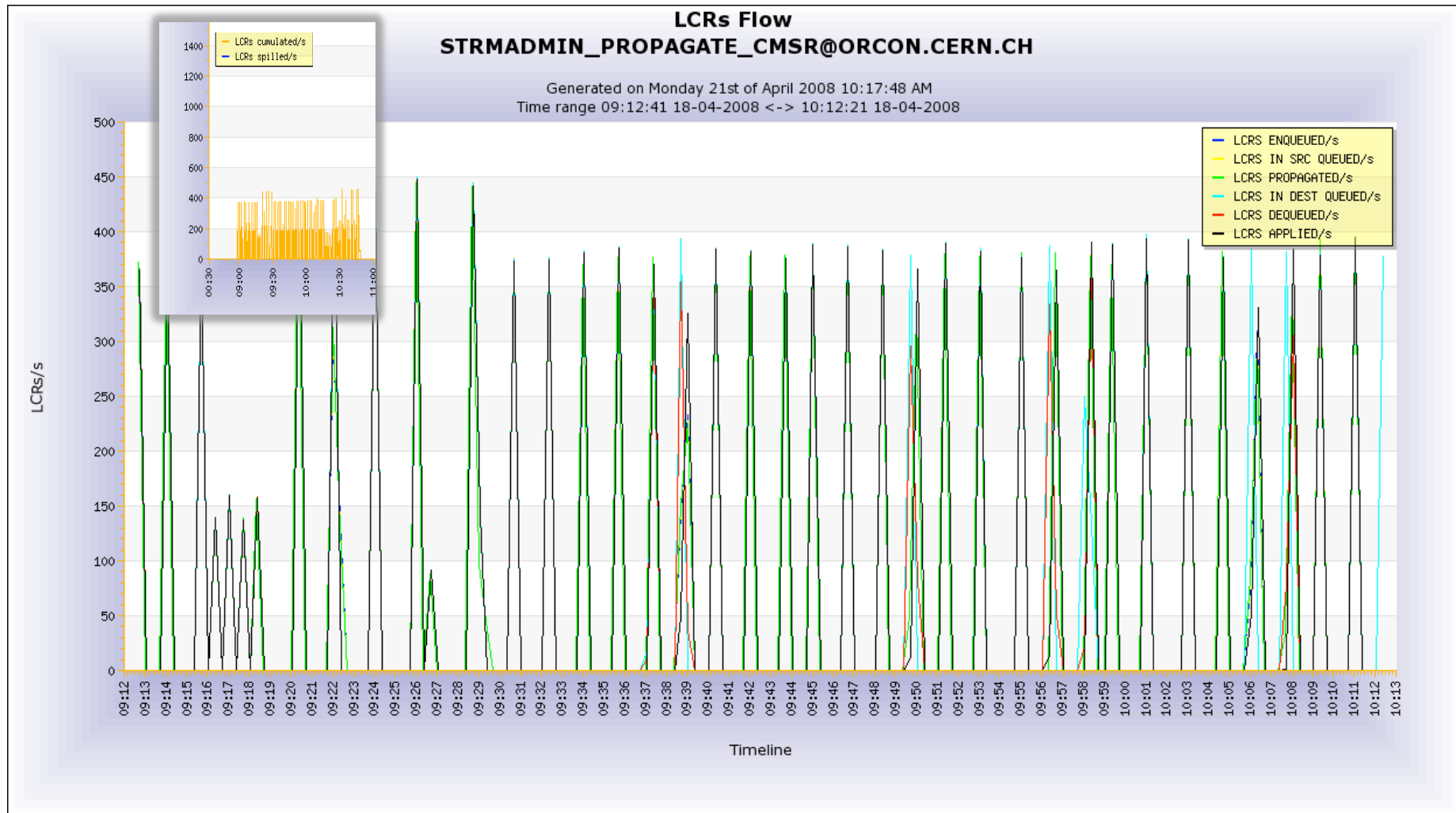
# Lessons Learned

- Example: CMS replication - Online to Offline (CERN)
  - single transaction mapping 428400 LCRs



# Lessons Learned

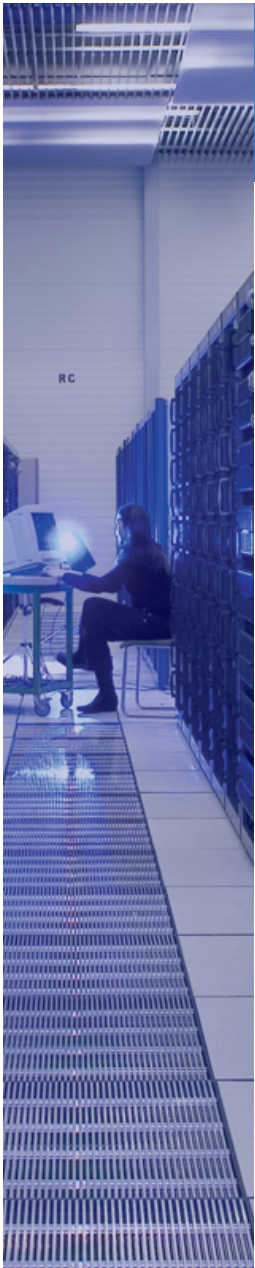
- Example: CMS replication - Online to Offline (CERN)
  - use BLOB objects: single transaction mapping 3600 LCRs

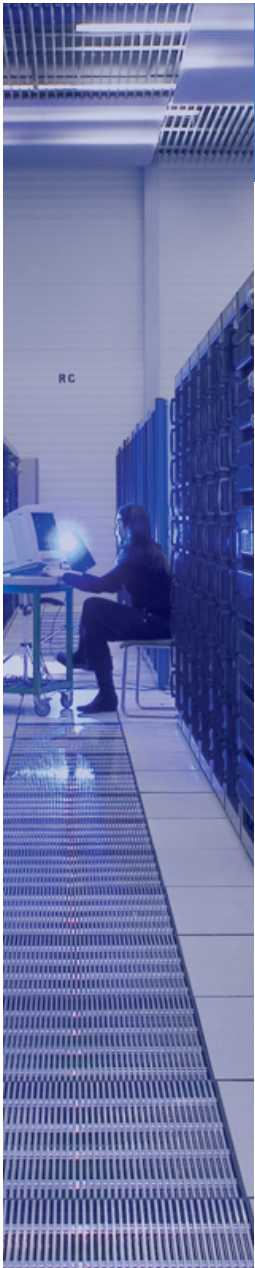




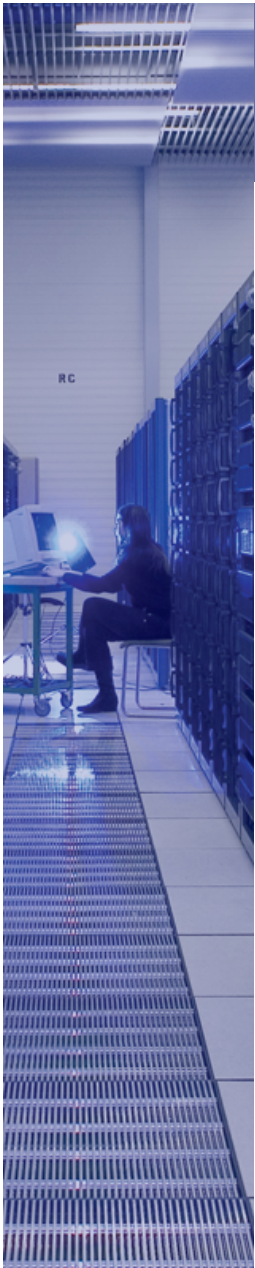
# Lessons Learned

- Apply oldest\_message\_number is not updated
  - caused by an old transaction not correctly removed from the apply spill table
  - dba\_apply\_spill\_txn view in order to identify the transaction
  - set the apply parameter `_IGNORE_TRANSACTION` with the transaction id in the apply spill over queue
  - run `purge_spill_txn` procedure (metalink note 556183.1)
- Apply might degrade performance when applying transactions to tables > 10M rows





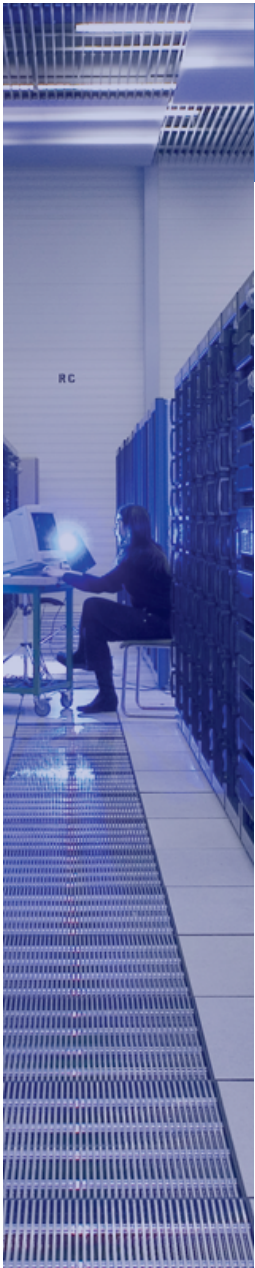
- How to **recover** Streams if downstream database crashes
  - use source database as replacement
    - all archive logs are available
  - check the oldest message number applied at each of the destinations
  - select Streams dictionary SCN < min(oldest message numbers)
  - create the Streams queue and all the propagations
  - create capture process where
    - first\_scn = dictionary SCN
    - start\_scn = oldest\_message\_number
- Configure **back** the downstream database
  - build a new Streams dictionary
  - stop capture and wait until all LCRs are applied
  - repeat steps above
  - register the archive logs with the capture process



- Performing a **switchover** in a Streams environment
  - database hw migration with minimal downtime
  - completely transparent for destination databases
  - source database:
    - before the switchover: move forward first\_scn
    - after the switchover: check that the archivelog files are registered with the capture process
      - otherwise, register them manually (from first\_scn)

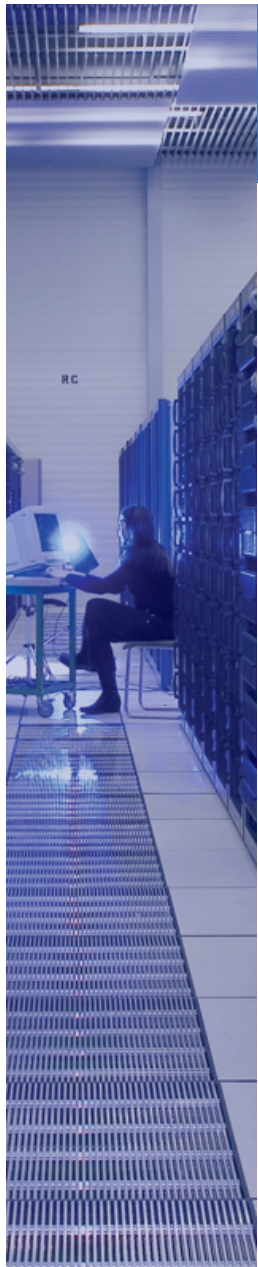
# Streams Bugs and Patches

- Streams specific patches [metalink note 437838.1](#)
- Bug 6452375
  - ORA-26687 in Streams from “drop table”
  - when two streams setups between same source and destination databases to replicate different schemas
- Bug 6402302
  - inconsistent capture/propagation/apply of DDLs in Streams
  - for example: “drop synonym” DDL is not captured/propagated or applied while create synonym is captured/propagated and applied



# Streams Bugs and Patches

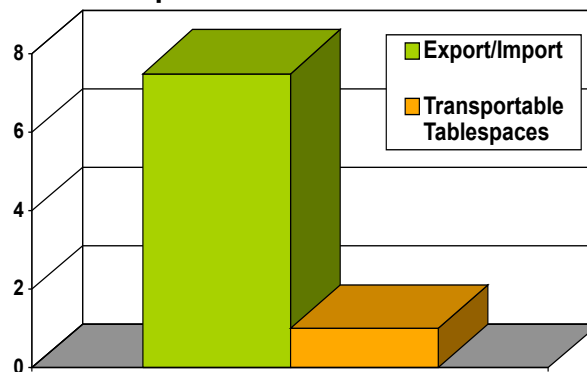
- ORA-00600: [KRVRDCBMDDLSQL1]
  - caused by rebuild index operation using parallel option
  - logminer corruption?
  - capture process could not be restarted at the current SCN
  - workaround proposed by Oracle: recreate capture using new dictionary after the index rebuild operation → data loss!!
  - complete re-instantiation of the Streams environment
- ORA-07445: exception encountered: core dump [kghufree()+485]
  - Oracle Database Change Notification cannot be used in a Streams environment



# Scalable Resynchronization

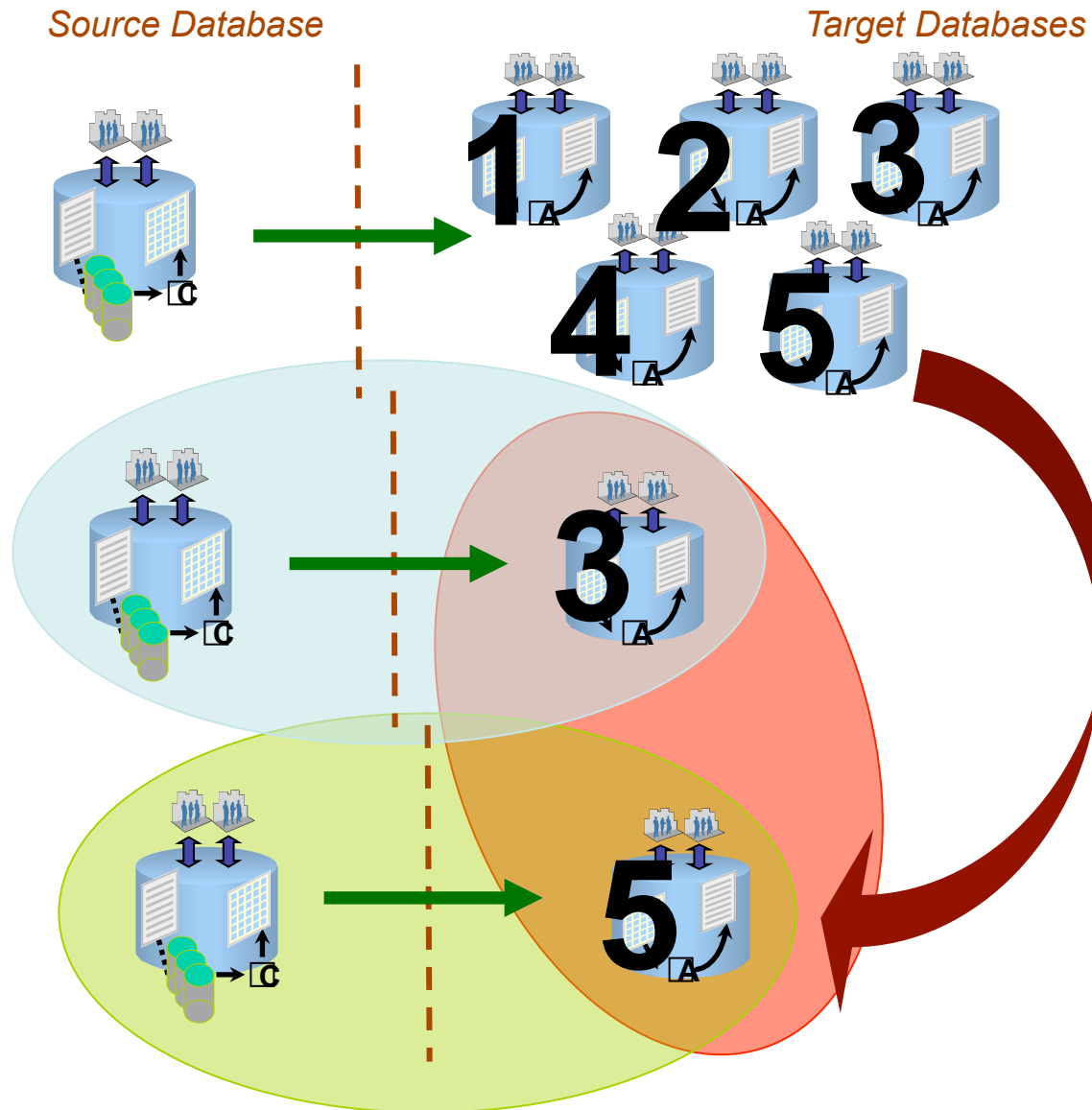
- Target site out of the Streams recovery window
- Complete transfer of data (schemas and tables) using Oracle Data Pump might take too long
  - Example ATLAS Conditions data
- **Transportable Tablespaces**: move a set of tablespaces from one Oracle database to another
  - Export metadata of tablespace instead of data in tablespace
- But tablespaces must be in **read-only** while the data is copied

Example: ATLAS COOL test data - 67 GB



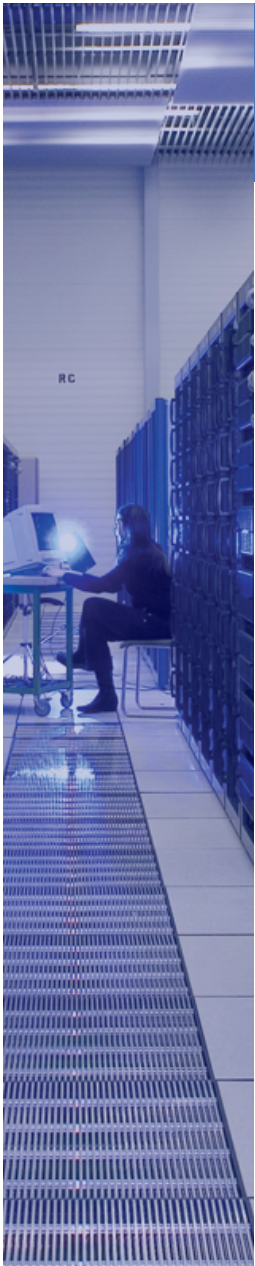
Moving data using transportable tablespaces is much faster than Data Pump export/import

# Scalable Resynchronization



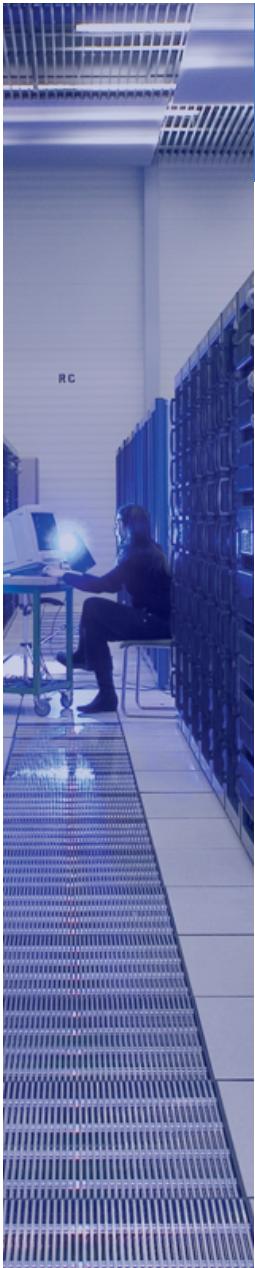
1. Create database links between databases
2. Create directories pointing to datafiles
3. Stop replication to site 5
4. Ensure tablespaces are read-only
5. Transfer the data files of each tablespace to the remote system
6. Import tablespaces metadata in the target
7. Make tablespaces read-write
8. Reconfigure Streams

- Oracle Enterprise Manager
  - Streams monitoring enhancements on 10.2.0.5
- Oracle Streams STRMMON monitoring utility
- Streams configuration report and health check script
  
- Extended tool for Streams monitoring: 3D Streams Monitor tool @CERN





# 3D Streams Monitor



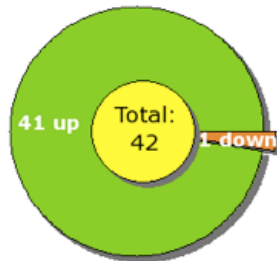
- Features:
  - Streams topology
  - Status of streams connections
  - Error notifications
  - Streams performance (latency, throughput, etc.)
  - Other resources related to the streams performance (streams pool memory, redo generation)
- Architecture:
  - “strmmon” daemon written in Python
  - End-user web application  
<http://oms3d.cern.ch:4889/streams/main>
- 3D monitoring and alerting integrated with WLCG procedures and tools

# 3D Streams Monitor

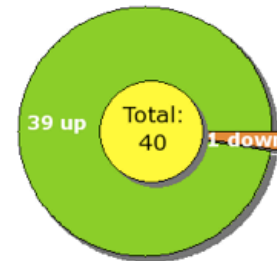
Monitor ▾ Maps ▾ DBs ▾ Streams ▾ Graphs ▾ Errors ▾ Availability ▾ History ▾ Reports ▾

## Monitor Summary

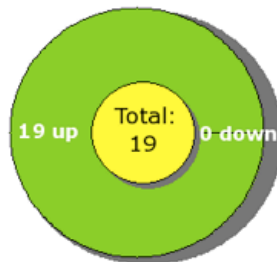
### Databases



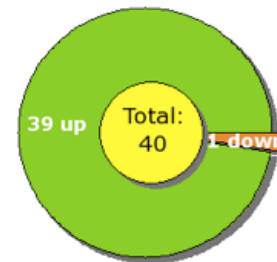
### Streams



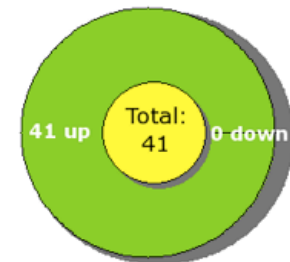
### Captures



### Propagations



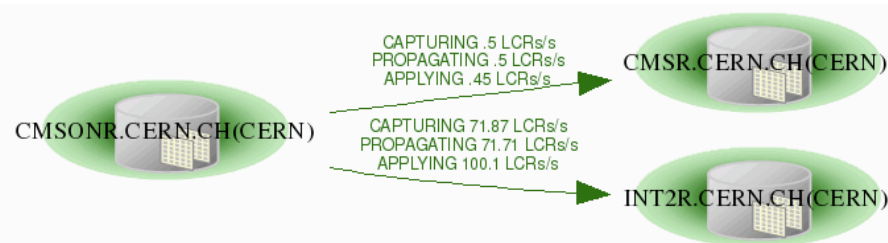
### Applies



Monitor ▾ **Maps** ▾ DBs ▾ Streams ▾ Graphs ▾ Errors ▾ Availability ▾ History ▾ Reports ▾

ALICE ▾ ATLAS ▾ **CMS** ▾ LHCb ▾ NEW ▾ TEST ▾ TEST11G2 ▾ TEST\_ATLAS ▾

TOPOLOGY auto refresh



# 3D Streams Monitor

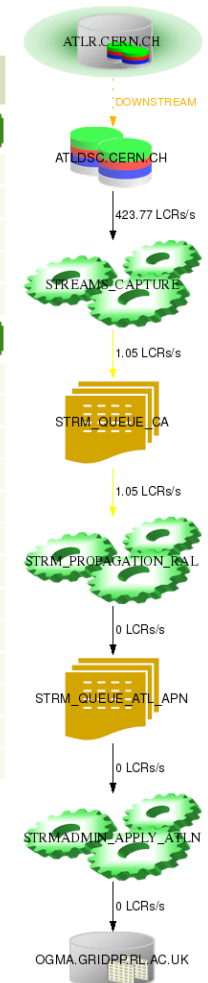
Monitor ▾ Maps ▾ DBs ▾ **Streams** ▾ Graphs ▾ Errors ▾ Availability ▾ History ▾ Reports ▾  
 All ▾ ALICE ▾ **ATLAS** ▾ CMS ▾ LHCb ▾ NEW ▾ TEST ▾ TEST11G2 ▾ TEST\_ATLAS ▾

## ACTIVE STREAMS

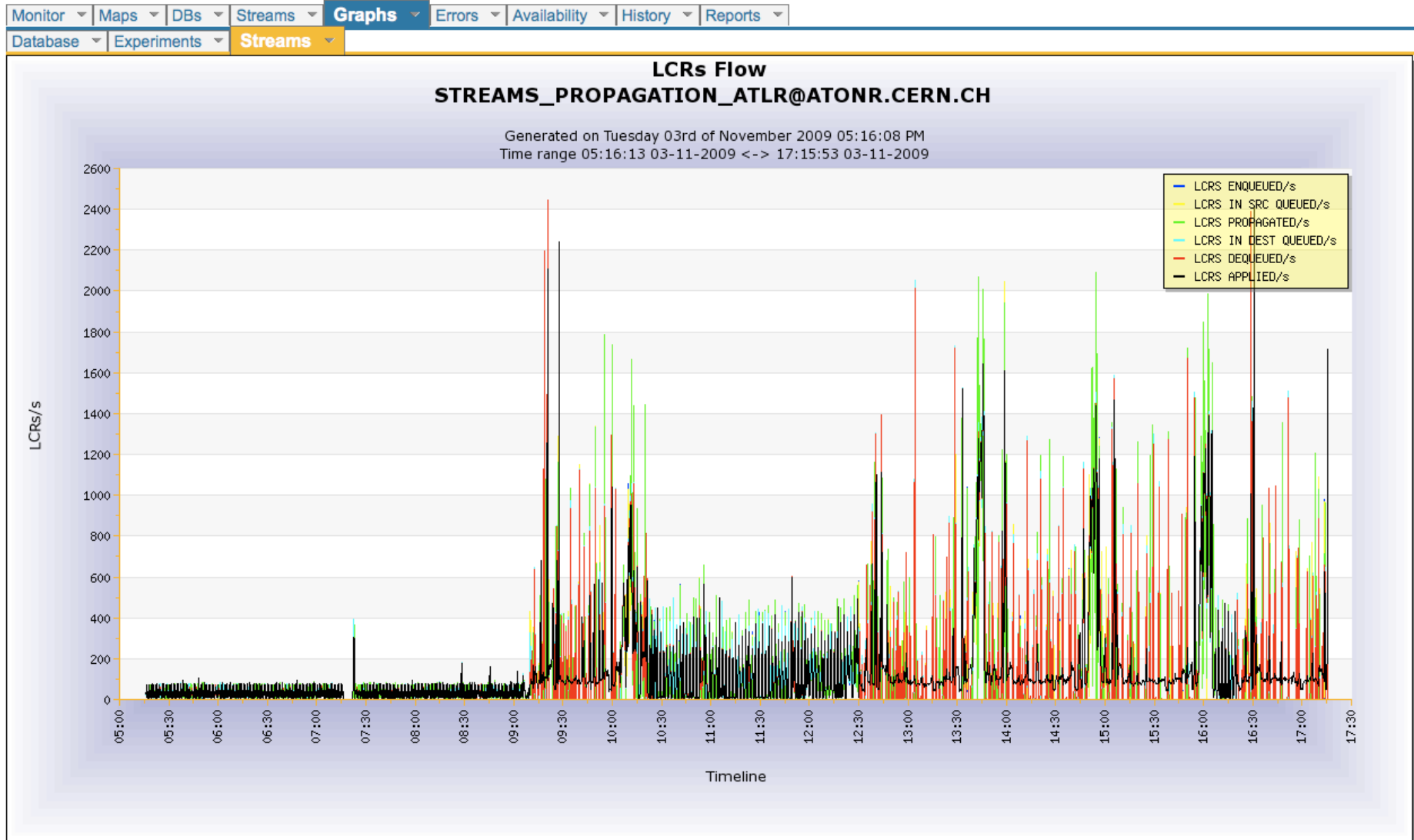
Stream	LCRs Cap	LCRs Enq	LCRs Prop	LCRs Deq	LCRs App	Cap Latency	Latency	Capture State	Propagation State	Apply State	State
<a href="#">ATLDSC.CERN.CH=&gt;ASGC3D.GRID.SINICA.EDU.TW</a>	951.34 /s	1.05 /s	1.05 /s	0 /s	0 /s	23 sec	--	CAPTURING CHANGES	ENABLED	IDLE	⚡

SOURCE		ATLDSC.CERN.CH
<b>CAPTURE</b>		
Name	STREAMS_CAPTURE	
SCN read	6083362640734	
LCRs captured	10491139 (542.77 /s)	
SCN captured	6083362640734	
LCRs enqueued	613528 (.25 /s)	
SCN enqueued	6083362621274	
Capture Latency	29 sec	
Info	ATLDSC.CERN.CH	
State	CAPTURING CHANGES	
Error Time		
Error Msg		
<b>QUEUE</b>		
Name	STRM_QUEUE_CA	
Id	42947	
Outstanding Msg	678227	
Cumulative Msg	4538140 (.25 /s)	
Cumulative Msg Spilled	711286 (0 /s)	
<b>PROPAGATION</b>		
Name	STRM_PROPAGATION_RAL	
LCRs Propagated	4532803 (.25 /s)	
Bytes Propagated	15453 (0 /s)	
State	ENABLED	
Error Time		
Error Msg		

DESTINATION		OGMA.GRIDPP.RL.AC.UK
<b>QUEUE</b>		
Name	STRM_QUEUE_ATL_APN	
Id	179048	
Outstanding Msg	0	
Cumulative Msg	94638683 (.05 /s)	
Cumulative Msg Spilled	158012 (0 /s)	
<b>APPLY</b>		
Name	STRMADMIN_APPLY_ATLN	
LCRs Dequeued	49704552 (.05 /s)	
Total LCRs Applied	48954298 (0 /s)	
SCN Dequeued	6083362647224	
Dequeue Latency	0 sec	
Transaction Received	528045	
Transaction Assigned	528045	
Transaction Applied	528045	
State	IDLE	
HWM SCN	6083362621274	
HWM Apply Latency	29 sec	
Error Time		
Error Msg		



# 3D Streams Monitor



# 3D Streams Monitor

Monitor | Maps | DBs | Streams | Graphs | Errors | Availability | History | **Reports**

## Reports Repository

### Streams report

Options

#### ALICE streams activities between 27.10.2009 and 03.11.2009

availability | instances | captures | propagations | applies | queues | undefined primary keys |

#### Capture processes stats

2009-10-27 <-> 2009-11-02

Legend:

LCRS\_CAP - LCRs captured  
LCRS\_ENQ - LCRs enqueued  
CAP\_LAT - capture latency - average time taken to capture single LCRs from log  
ENQ\_LAT - enqueue latency - average time taken to enqueue single LCRs (capture time included)

max min

2009-10-27				2009-10-28				2009-10-29				2009-10-30				2009-10-31				2009-11-01				2009-11-02					
DB_NAME	PROCESS_NAME	LCRS_CAP	LCRS_ENQ	CAP_LAT	ENQ_LAT	LCRS_CAP	LCRS_ENQ	CAP_LAT	ENQ_LAT	LCRS_CAP	LCRS_ENQ	CAP_LAT	ENQ_LAT	LCRS_CAP	LCRS_ENQ	CAP_LAT	ENQ_LAT	LCRS_CAP	LCRS_ENQ	CAP_LAT	ENQ_LAT	LCRS_CAP	LCRS_ENQ	CAP_LAT	ENQ_LAT				
ALIONH.CERN.CH	STREAMS_PVSS_CAPTURE	8766353	291876	5.5 sec	3.6 sec	7418204	241186	5.4 sec	3.7 sec	6742596	599314	5.2 sec	3.5 sec	6799907	186108	5.3 sec	3.7 sec	6317804	193119	5.3 sec	3.7 sec	6983790	201615	5.2 sec	3.7 sec	7591331	204735	5.4 sec	3.6 sec

Total

TOTAL					
DB_NAME	PROCESS_NAME	LCRS_CAP	LCRS_ENQ	CAP_LAT	ENQ_LAT
ALIONH.CERN.CH	STREAMS_PVSS_CAPTURE	50215988	29759103	5.3 sec	2d 11h 2m 45s

#### Propagation processes stats

2009-10-27 <-> 2009-11-02

Legend:

LCRS\_PROP - number of propagated LCRs  
BYTES\_PROP - number of propagated bytes

max min

2009-10-27			2009-10-28			2009-10-29			2009-10-30			2009-10-31			2009-11-01			2009-11-02		
DB_NAME	PROCESS_NAME	LCRS_PROP	BYTES_PROP	LCRS_PROP	BYTES_PROP	LCRS_PROP	BYTES_PROP	LCRS_PROP	BYTES_PROP	LCRS_PROP	BYTES_PROP	LCRS_PROP	BYTES_PROP	LCRS_PROP	BYTES_PROP	LCRS_PROP	BYTES_PROP			
ALIONH.CERN.CH	STREAMS_PROPAGATE_POBR	5827071	130.3 kB	4814418	80.2 kB	3179169	79.7 kB	3714992	130.4 kB	3854643	105.3 kB	4024718	0 B	4086873	92.0 kB					

Total

TOTAL			
DB_NAME	PROCESS_NAME	LCRS_PROP	BYTES_PROP
ALIONH.CERN.CH	STREAMS_PROPAGATE_POBR	29701885	617.9 kB

#### Queue stats

2009-10-27 <-> 2009-11-02

Legend:

LCRS\_IN\_Q - number of enqueued messages (LCRs)  
LCRS\_SPILLED - number of spilled LCRs to disk  
SPILL\_RATE - spilling percentage  
LCRS\_AVG - average number of messages in queue

max min

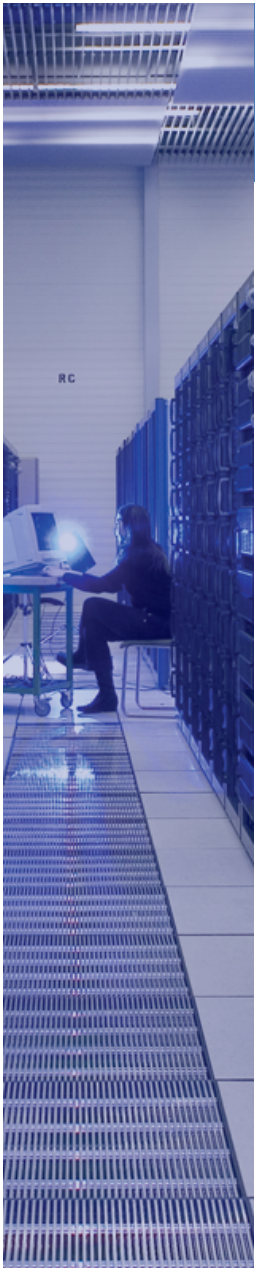
2009-10-27				2009-10-28				2009-10-29				2009-10-30				2009-10-31				2009-11-01				2009-11-02					
DB_NAME	QUEUE_NAME	LCRS_IN_Q	SPILLED	SPILL_RATE	LCRS_AVG	LCRS_IN_Q	SPILLED	SPILL_RATE	LCRS_AVG	LCRS_IN_Q	SPILLED	SPILL_RATE	LCRS_AVG	LCRS_IN_Q	SPILLED	SPILL_RATE	LCRS_AVG	LCRS_IN_Q	SPILLED	SPILL_RATE	LCRS_AVG	LCRS_IN_Q	SPILLED	SPILL_RATE	LCRS_AVG				
ALIONH.CERN.CH	STREAMS_PVSS_QUEUE_CA	5834670	0	0.0 %	660	4821164	0	0.0 %	560	3384564	0	0.0 %	397	3720078	0	0.0 %	415	3860081	0	0.0 %	405	4030247	0	0.0 %	450	4002501	0	0.0 %	471
POBR.CERN.CH	STREAMS_PVSS_QUEUE_AP	5866178	0	0.0 %	135	4840116	0	0.0 %	144	3396840	0	0.0 %	96	3737837	0	0.0 %	106	3876122	0	0.0 %	81	4046299	0	0.0 %	100	4110115	0	0.0 %	111
POBR.CERN.CH	KUPCBS_1_20091015143048	0	0	0.0 %	0	0	0	0.0 %	0	0	0	0.0 %	0	0	0	0.0 %	0	0	0	0.0 %	0	0	0	0.0 %	0	0	0	0.0 %	0

Total

TOTAL					
DB_NAME	QUEUE_NAME	LCRS_IN_Q	SPILLED	SPILL_RATE	LCRS_AVG
ALIONH.CERN.CH	STREAMS_PVSS_QUEUE_CA	29743307	0	0.0 %	480.2
POBR.CERN.CH	STREAMS_PVSS_QUEUE_AP	29873509	0	0.0 %	110.8
POBR.CERN.CH	KUPCBS_1_20091015143048	0	0	0.0 %	0.0

Generated in 39 min 57 sec

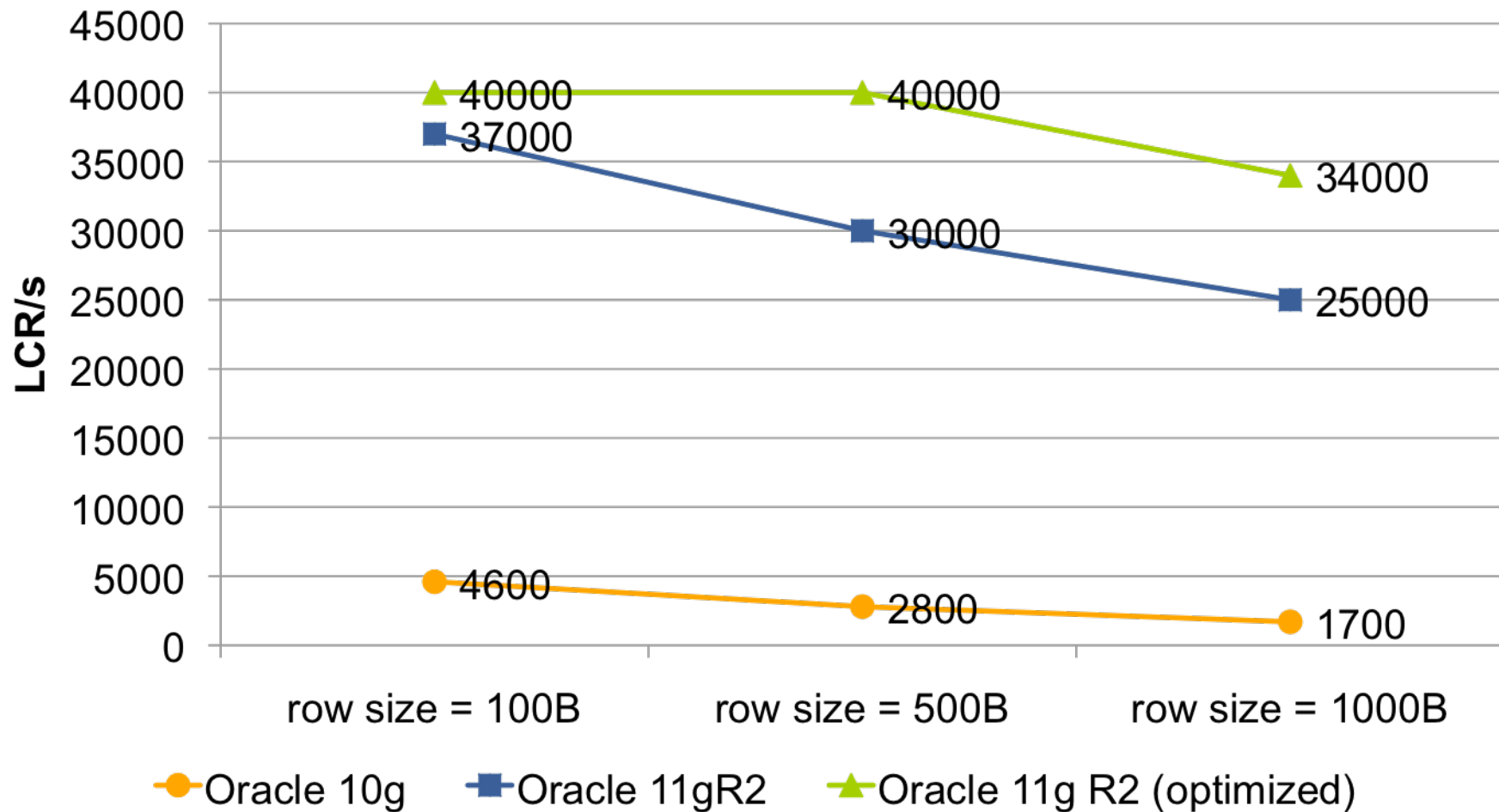
# New Streams 11g Features



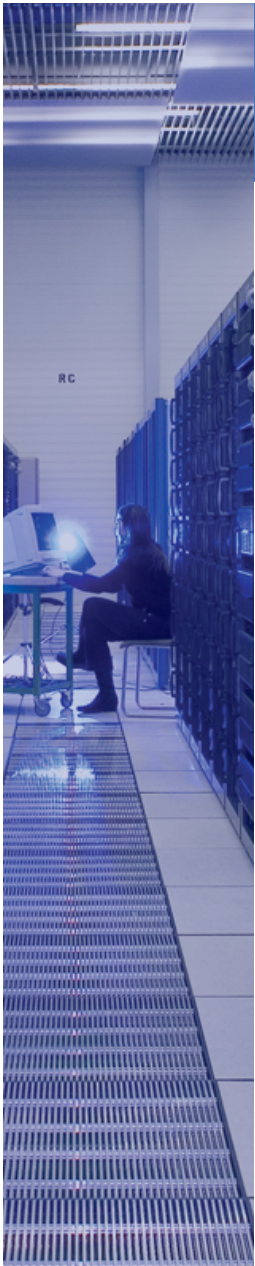
- Performance improvements:
  - reader process mines from the in-memory redo log buffers
    - minimizes disk I/O
    - reduces capture latency
  - direct communication between capture and apply processes: **Combined Capture and Apply**
    - improves LCR transmission throughput
    - reduces end-to-end replication latency
  - internal mechanism to execute change records and extensive caching
    - reduces CPU consumption
    - minimizes latch contention and other wait events

# New Streams 11g Features

## Average Streams Throughput



# New 11g Streams Features



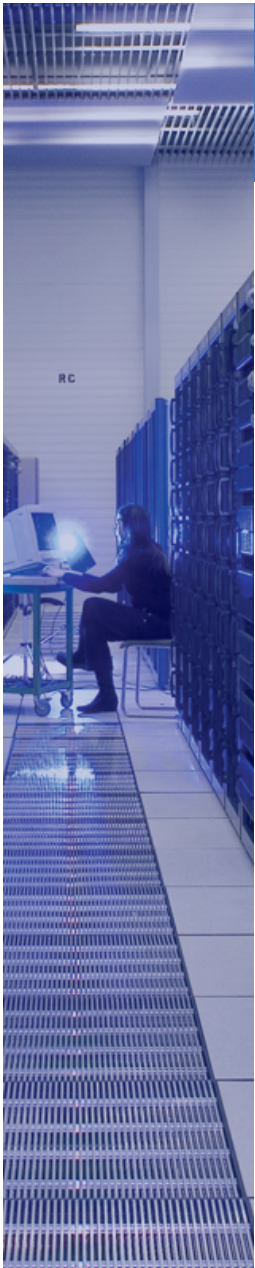
- Automatic Split and Merge
  - split a stream in cases where a replica is unavailable
  - merge into a single stream when replica catches up
  - procedures and sql script generation
  - automatic replication management based on thresholds
- Compare and Converge
  - compare objects across databases for inconsistency
  - resynchronize objects if required
  - table or column level synchronization
  - additional scripting for schema comparison



# New 11g Streams Features

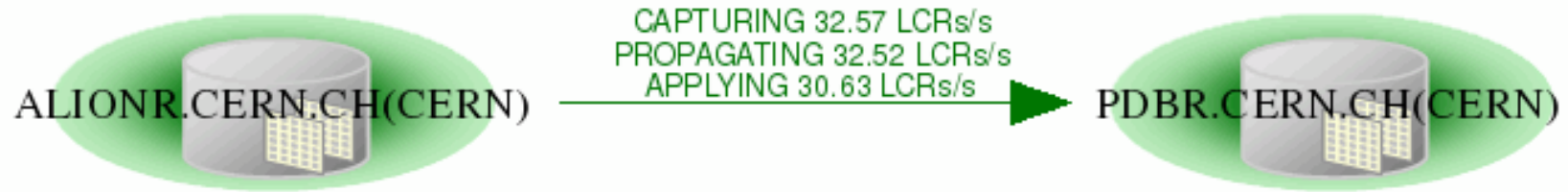
And more...

- Synchronous capture
- LCRs track through a Stream
- Topology and Performance Advisor
- New error messages for error handling

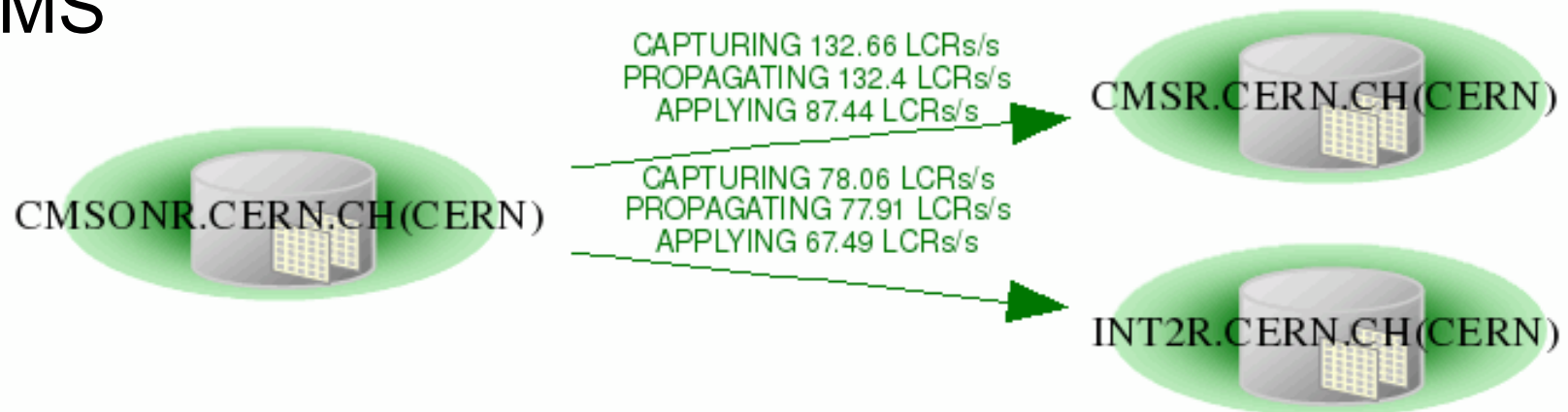


# More Streams Setup Examples

- ALICE

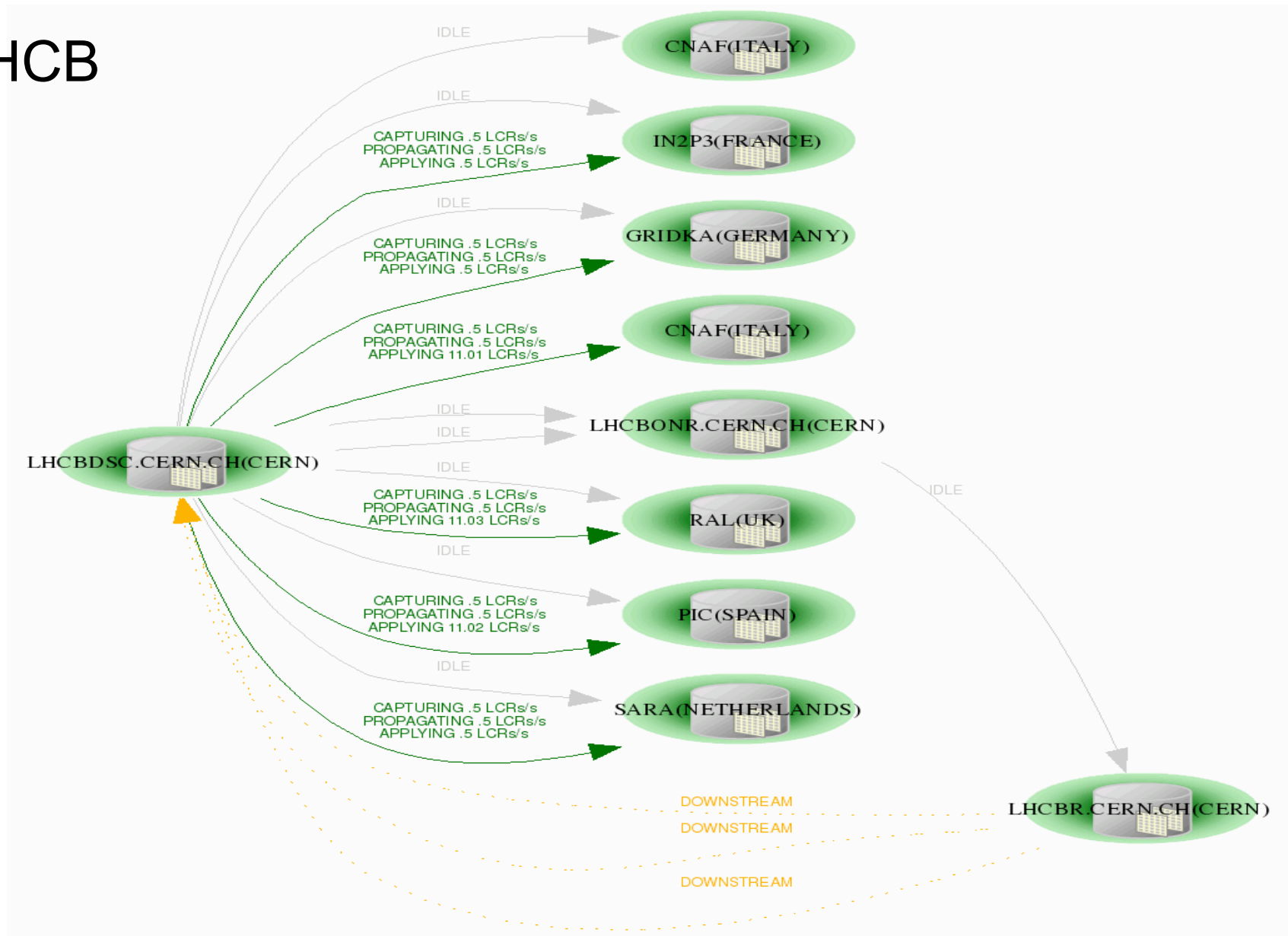


- CMS

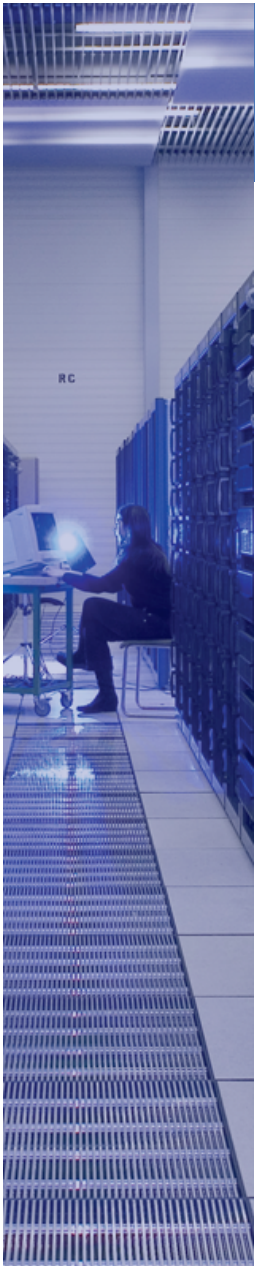


# More Streams Setup Examples

- LHCB



# Summary



- The LCG Database Deployment Project (LCG 3D) has set up a world-wide distributed database infrastructure for LHC
  - some 33 RAC clusters = 636 CPU cores at CERN + several tens of nodes at 10 partner sites are in production now
- Large scale tests have validated that the experiment are implemented by the RAC & streams based set-up
  - backup & recovery tests have been performed to validate the operational procedures at all sites
- Monitoring of database & streams performance has been implemented building on grid control and strmmmon tools
  - key to maintain and optimize any larger system

# Q&A

