

Docker @ CDS

André Schaaff¹, François-Xavier Pineau¹, Gilles Landais¹, Laurent Michel²

¹Centre de Données astronomiques de Strasbourg, ²SSC-XMM-Newton

Paul Trehou

Université de technologie de Belfort-Montbéliard

IVOA, Shanghai, 14-19/05/2017



H2020-Astronomy ESFRI and Research Infrastructure Cluster (Grant Agreement number: 653477).

□ Docker for VizieR deployment

- VizieR is deployed on several mirrors
- Hosts with different Linux distributions, kernels, etc.
- Docker as a solution to deploy “quick and easy” ?

The screenshot displays the VizieR web interface. At the top, there is a navigation bar with the CDS logo and links for Portal, Simbad, VizieR, Aladin, and X-M. Below the navigation bar, the main content area is divided into several sections. On the left, there is a 'Search Criteria' section with a 'Preferences' sub-section containing a 'max:' dropdown set to 50, an 'HTML Table' dropdown, and an 'All columns' checkbox. Below this is a 'Mirrors' section with a dropdown menu currently showing 'CDS, France' and a list of other mirrors including Tokyo, Japan; CADAC, Canada; Cambridge, UK; CFA/Harvard, USA; UKIRT-Hawaii, USA; IUCAA, India; Beijing Obs., China; and SAAO., South Africa. On the right, there is a search area titled 'Find catalogs among 1580' with a 'Clear' button and a search input field. Below this is a section for 'Search by Position across' with a 'Target Name (resolved by Sesar)' field and another 'Clear' button. At the bottom right, there is a link for 'More about VizieR' and a section titled 'Tools related to VizieR'.

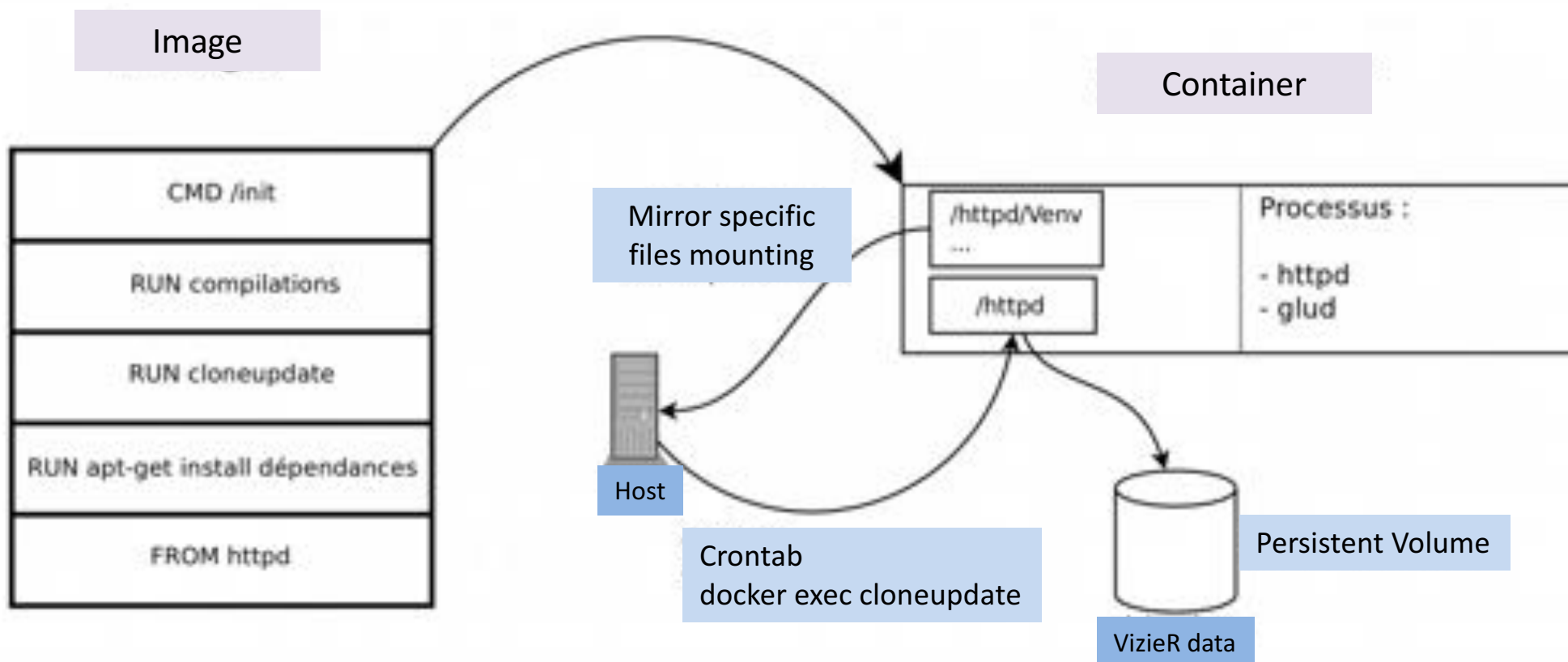
□ VizieR mirror update process

- Typical installation of VizieR
 - Dependencies and CGI scripts transfer through Rsync and scp
 - Dependencies installation with the package manager
 - Compilation of the dependencies only available as sources (developed at CDS)
 - CGI files copy and Apache configuration
 - Apache start

□ To a VizieR Docker image...

- Prototyping following the previous process
 - Resolution of missing packages (like gcc, make, rsync, .., not present in the Apache image)
- To an optimized version
 - Sources directly from CVS
 - CGI scripts and static files in a Docker volume (between the host and the container), these files will be updated through scripts (executed at each container start and via a cronjob service (on the host))
 - Path corrections
 - Also use of Portainer (Simple management UI)

□ VizieR & Docker at the End



□ And now

- Used in a first step inside CDS
 - To test Docker use on a mid-term basis
 - To install local “VizieR” to test it before update
- Needs discussion and agreement with at least one host (in a first step) to start to use it “in production” all over the world !

Experiments with Spark - status

André Schaaff, François-Xavier Pineau
Centre de Données astronomiques de Strasbourg

Osman Aidel
Centre de Calcul – IN2P3 Villeurbanne

Noémie Wali, Paul Trehou
Université de technologie de Belfort-Montbéliard

IVOA, Shanghai, 14-19/05/2017



H2020-Astronomy ESFRI and Research Infrastructure Cluster (Grant Agreement number: 653477).

□ Outline

Apache Spark

Use case & data

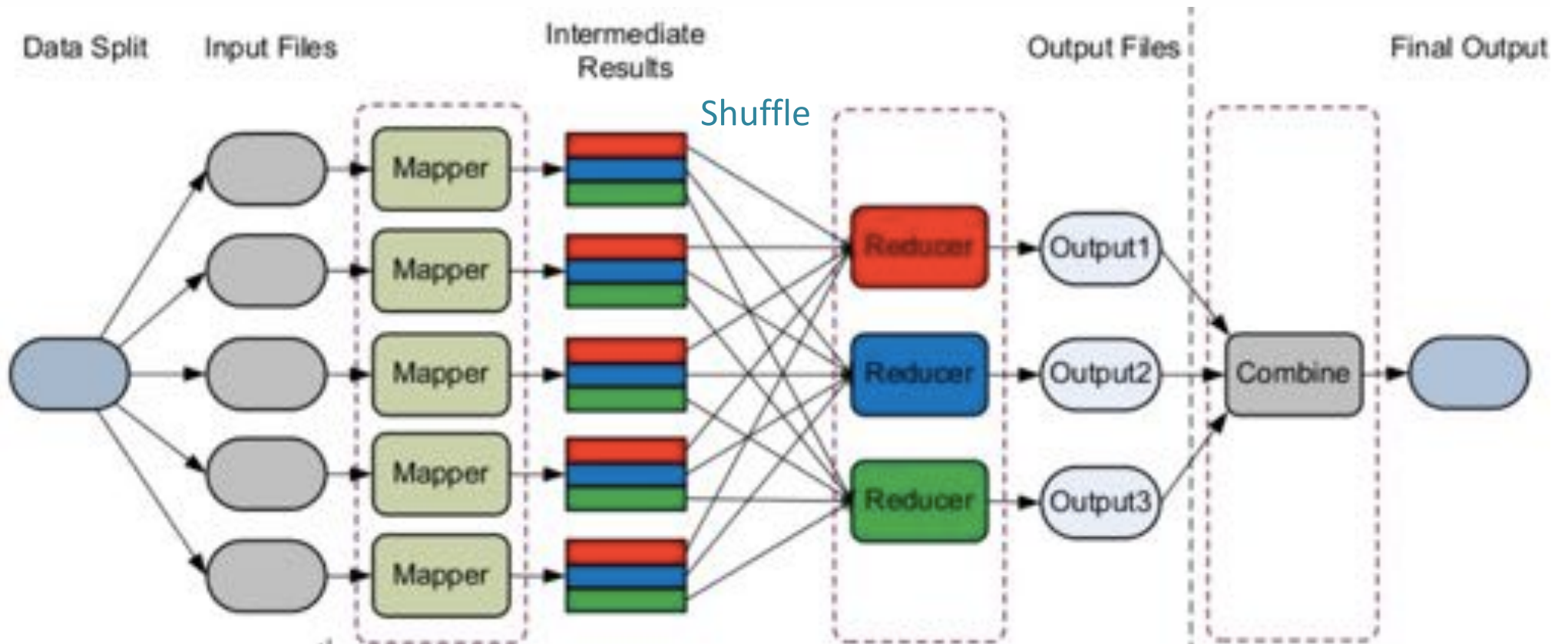
Test beds, experiments and what we learned

Perspectives

□ Apache Spark

- “Apache Spark is a cluster computing platform designed to be fast and general purpose.”
- It extends the MapReduce model to support more types of computations (interactive queries, stream processing, etc.) and it offers APIs for Scala, Java, Python, R,...

□ MapReduce



Credit: G. Fedak, INRIA

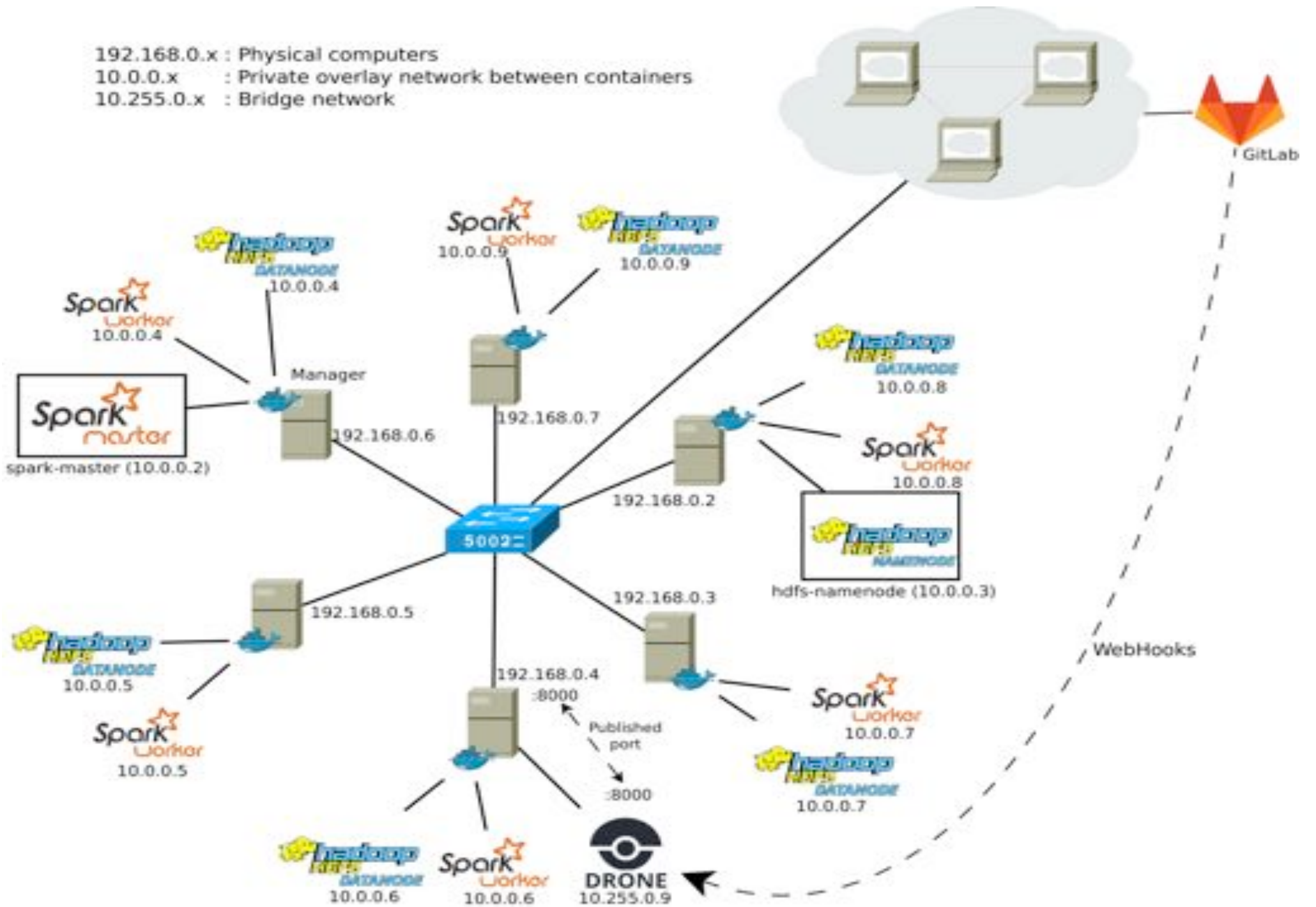
□ Apache Spark, so quick ?

- **Computations in memory** (as much as possible, otherwise pilling to the disks)
- Introduction of data models
 - **RDD** (Resilient Distributed Datasets)
 - **Immutable** distributed collection of elements
 - Operations: **Transformations** (map, filter, etc.), **Actions** (reduce, count, etc.)
 - **Datasets** to represent tabular data, queryable via SQL
- It uses mainly Hadoop Distributed File System (HDFS).

□ Other technical aspects

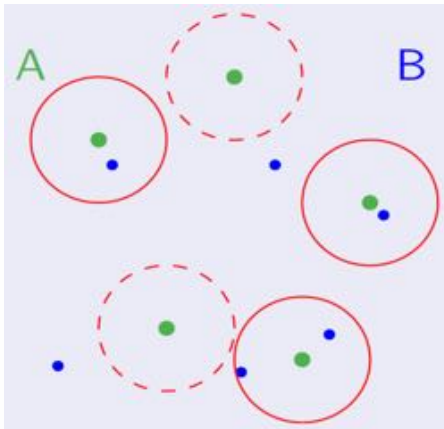
- Introduction of **Docker** (components) and **Drone** (continuous integration) to “**automate**” the deployment process and to focus mainly on the development side. It is becoming easy to migrate to external resources when needed.
- Use of **Scala** which is native in Spark (a part of the Java API is “experimental”).

192.168.0.x : Physical computers
 10.0.0.x : Private overlay network between containers
 10.255.0.x : Bridge network



□ Use case & data

- The “cross-match” of (large) source catalogues.
- Examples:
 - 2MASS¹, 470,992,970
 - SDSS² DR9, 469,053,874



Full sky: all the sources
 A cone: only the sources which are at a certain angular distance from a given position
 A HEALPix cell



Fuzzy join between 2 tables (A and B) of several hundred millions of data

Byte	Format	Units	Label	Explanations
1-	10	F10.6	deg	RAdeg [ra] Right ascension (J2000)
12-	21	F10.6	deg	DEdeg [dec] Declination (J2000)
23-	26	F4.2	arcsec	errMaj [err_maj] Semi-major axis of error ellipse
28-	31	F4.2	arcsec	errMin [err_min] Semi-minor axis of position error ellipse
33-	35	I3	deg	errPA [err_ang] Position angle of error ellipse [0,180]

VizieR Result Page

The 3 columns in color are computed by VizieR, and are not part of the original data.

1/246out 2MASS All-Sky Catalog of Point Sources (Cutri, 2003) 2013CaJ204...0C ResNetStu

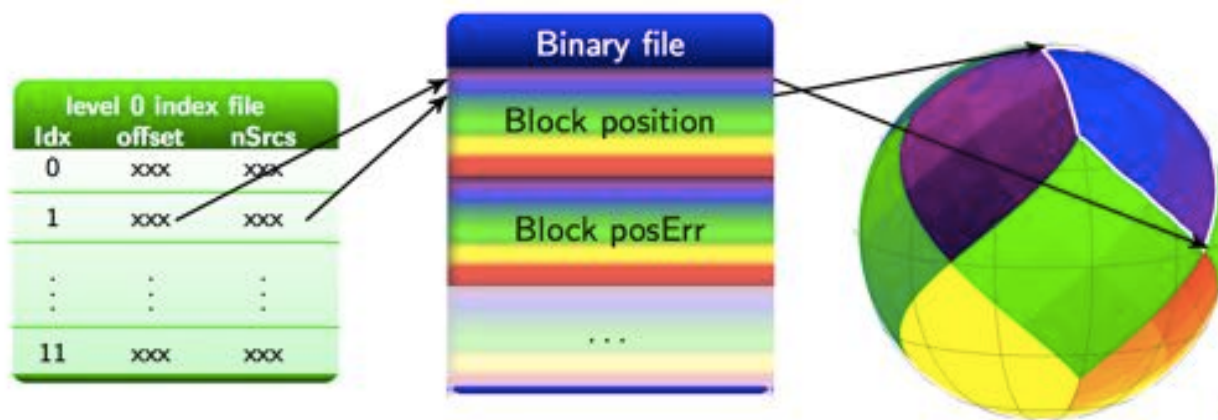
2MASS All-Sky Catalog of 470,992,970 sources. Please acknowledge the usage of the 2MASS All-Sky Survey; see also the 2MASS Pages. Note that the magnitudes in red correspond to low quality results (upper limits or very poor photometry) (470992970 rows)

RAJ2000	DECJ2000	RAJ2000	DECJ2000	2MASS	Jmag	Hmag	Kmag	Q16	H16	B16	C16	X16	A16	
210.0171100	-42.44.3317	+41.16.08.53	010.684737	+41.269035	00424433+4116085	9.433	0.052	8.568	0.051	8.473	0.051	8.222	111.000	2.0

Credits: <http://healpix.jpl.nasa.gov/>

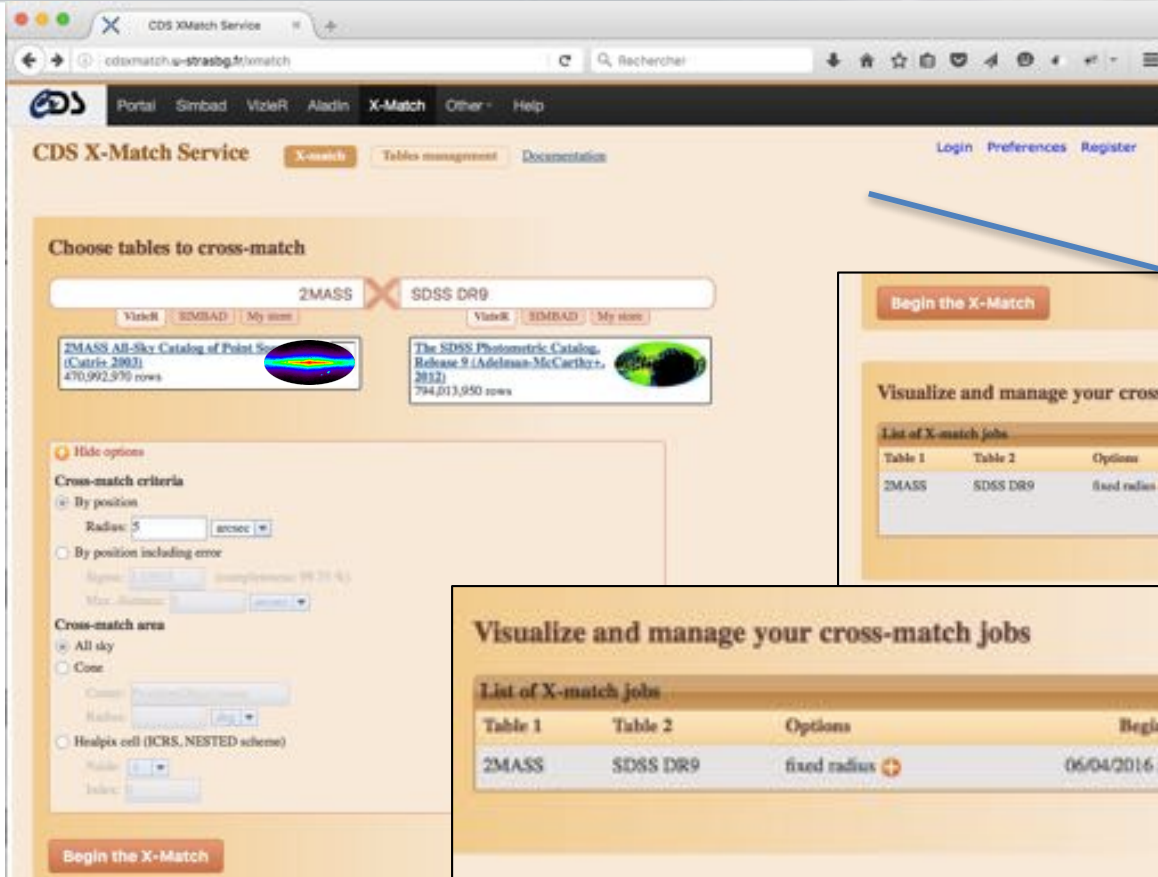
□ Not distributed but...

- organised and stored on one server (2x10 cores, 64GB, 12TB (15k tours))



The sky is cut into diamonds of the same size, pixels, each source or sky object is a numbered pixel.

Illustration: X-Match frontend



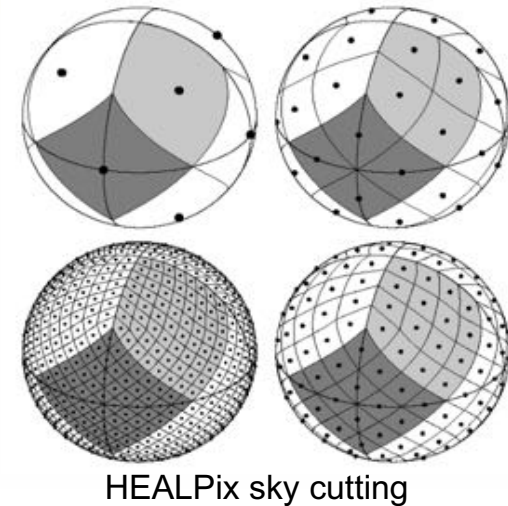
X-Match of
2MASS & SDSS DR9
(over 10,000 catalogues + own
catalogue upload)



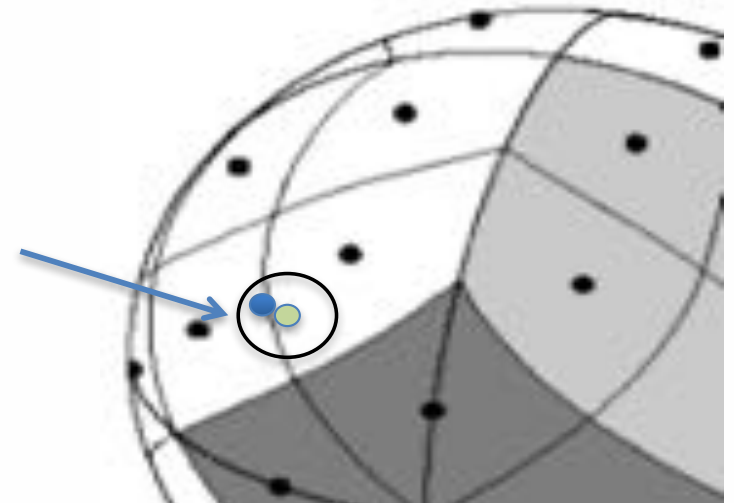
GAIA DR1 X SDSS DR9 (1 arcsec) in 17' (100.10⁶ matches, 34 GB)

□ Illustration

- A X-Match implementation in MapReduce, couples (Key = pixel number, Value)



- Side effects
 - Fuzzy join
 - Source duplication in the neighbour cells if needed



Credits: HEALPix – arXiv:astro-ph/0409513

□ Test beds: hardware & software

- **External resources** (Experiment 1)
 - 12 nodes, 4 cores, 32GB, Raid 2*2TB, Ubuntu
 - Configuration was defined “ad hoc” and low cost
- **Collaboration with IN2P3** (Experiment 2)
 - 9 nodes (only one VM per physical server, CentOS / OpenStack, remote data storage), 24 threads (-> 216 threads), 64GB (-> 576 GB) -> possible X-Match of **billion sources**
- **Internal resources** to prototype, (Experiment 3)
 - 6 nodes (4 cores, 16GB, 1TB + 500GB **SSD**), Ubuntu
- Apache distributions of **Spark and Hadoop, Java, Scala, Docker, Drone, ...**

□ Experiment 1 (12 nodes)

- Input data ([SDSS DR7](#) (primary sources) and [2MASS](#)): 54GB and 58GB file size; 357 175 411 and 470 992 970 elements
- Output data: 49 208 820 elements

X-Match service reference time was: 10 minutes

Cross-Match (source duplication done in phase 2 with all the data as output)					
HDFS block size= 128MB for the input files ; sdss7.csv and t 2mass.csv replicated 2 times					
HashPartitioner	60 partitions				
HDFS output files size	32MB				
Number of nodes Spark/HDFS	5	7	9	10	11
Phase 1: prepare	23,0	16,0	14,0	14,0	13,0
mapToPair (sdss7.csv)	5,1	4,9	4,9	4,8	4,7
saveAsHadoopFile (sdss7.bin)	5,7	2,7	2,0	2,3	1,5
mapToPair (2mass.csv)	5,7	5,2	5,2	5,1	5,0
saveAsHadoopFile (2mass.bin)	6,5	3,6	1,9	1,6	1,4
Phase 2: join	31,0	21,0	13,0	11,0	9,9
mapToPair (sdss7.bin)	7,2	4,7	3,5	3,0	2,6
flatMapToPair (2mass.bin)	11,8	8,3	5,5	4,9	4,3
saveAsTextFile (crossMatch_D.txt)	12,0	7,6	3,4	2,4	2,3
TOTAL	54,0	37,0	27,0	25,0	22,9

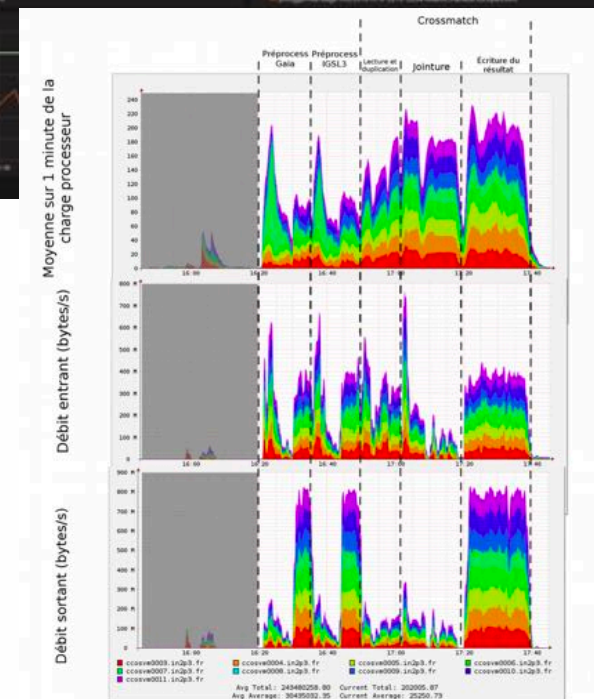
□ What we have learned

- Time was similar to the X-Match service from 11 nodes
- Keys common to 2 RDDs are not necessarily on the same node
 - It implies a **transfer overhead** between the nodes during the join => impact on the **performances**
 - We had clearly a **bottleneck** in the join phase (“**shuffle**”)
 - “block affinity groups” is an on-going work at Apache.
- We spent time on the “data co-location”
- We found a solution to do it “**manually**” via Python scripts (=> Experiment 3, with SSD).

□ Experiment 2: IN2P3 cluster

- Gaia (1.1 billion sources) X IGSL3 (1.2 billions)
- 1.6 billion associations in 15 minutes (30 minutes for the production X-Match Server)*
- Debriefing: Bottleneck is
 - Network for Loading phase
 - Mainly CPU during the xmatch phase (can be improved)

* Rebuild Join not included



□ Experiment 3: SSD

- Network 1GB/s, SSD 400-450 MB/s
- We were able to experiment the Python Script to move the data blocks « manually » (=> manual colocation of data)
 - With the HD, no gain (Network and HDD have a similar speed)
 - With the SSD, a gain of 30% on the loading phase, as theoretically predicted (SSD faster than network)
 - And similar performances in or out of Docker

□ Perspective and conclusion

- Apache Spark quick to install, easy to use for common tasks
- But not easy to understand what happens, how it works when you have particular use cases
- Real interest in several communities
- Ongoing experiments, probably with clouds
- Docker use will continue with an application to the X-Match service

□ Links

- Apache Spark, <http://spark.apache.org/>
- Apache Hadoop, <http://hadoop.apache.org/>
- Spark : Cluster Computing with Working Sets, Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, Ion Stoica, University of California, Berkeley, http://static.usenix.org/legacy/events/hotcloud10/tech/full_papers/Zaharia.pdf
- Optimizing Shuffle Performance in Spark, Aaron Davidson, Andrew Or, UC Berkeley, http://www.cs.berkeley.edu/~kubitron/courses/cs262a-F13/projects/reports/project16_report.pdf
- Resilient Distributed Datasets : A Fault-Tolerant Abstraction for In-Memory Cluster Computing, Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica, University of California, Berkeley, https://www.cs.berkeley.edu/~matei/papers/2012/nsdi_spark.pdf
- JavaSpark Api, <http://spark.apache.org/docs/latest/api/java/>
- HEALPix, <http://healpix.jpl.nasa.gov/>



H2020-Astronomy ESFRI and Research Infrastructure Cluster (Grant Agreement number: 653477).