# Mining Archives: needs for Machine Learning

S. Cavuoti & G. Riccio

INAF - Astronomical Observatory of Capodimonte

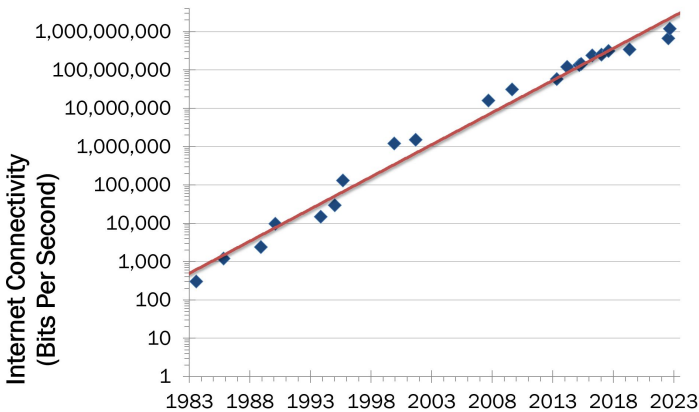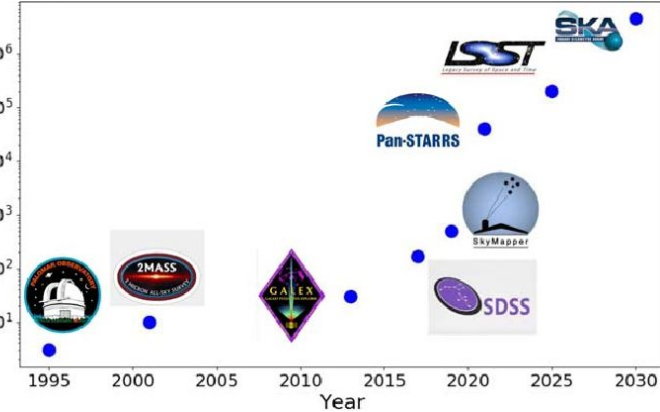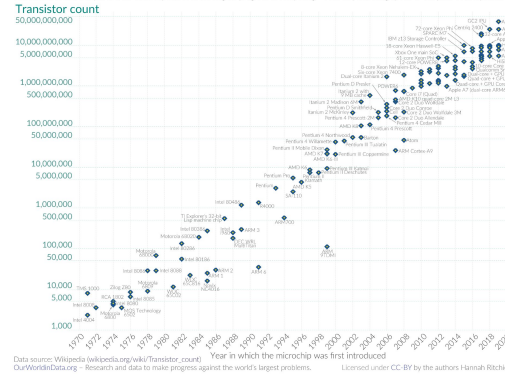# Users do what users do best

# Simple Premise



There are archives already implementing
some (or even all) solutions to the problems
that I will address, there is no need to jump up
and say: "*I'm not guilty*"

# Data Tsunami

We are all facing Data Tsunami

# Data Tsunami
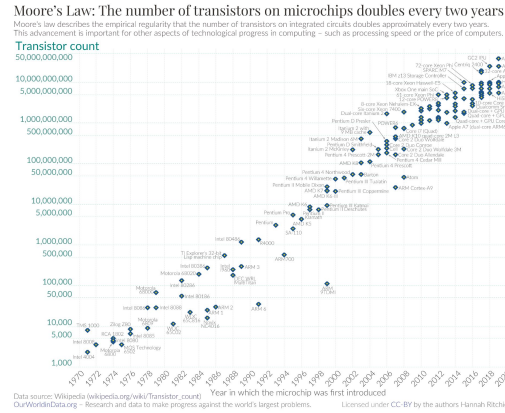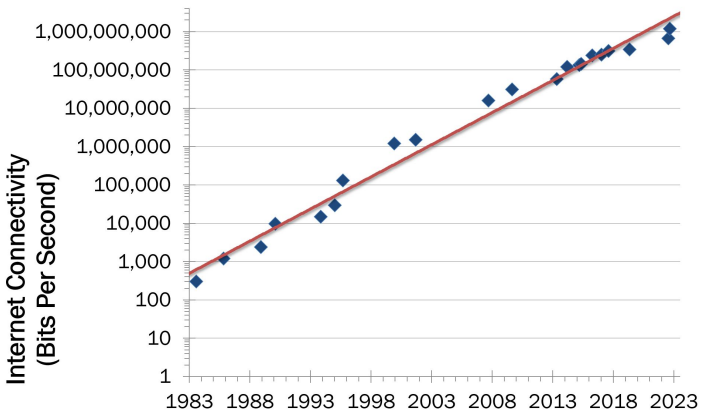
We are all facing Data Tsunami





Moore's Law: The number of transistors on microchips doubles every two years
Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years.
This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.

Transistor count

Year in which the microchip was first introduced

Data source: Wikipedia (wikipedia.org/wiki/Transistor_count)
OurWorldinData.org – Research and data to make progress against the world's largest problems.       Licensed under CC-BY by the authors Hannah Ritchie

# Data Tsunami

Astronomers are increasing too, even though not at the same speed



**More users requiring data, that are by themself increasing faster than the connection.**

# Data Tsunami

Transferring large datasets from astronomical archives to local machines is inefficient, consuming excessive bandwidth and storage.

One possible solution could be instead of moving data to the code, move the code to the data by providing computational resources directly within the archive infrastructure.

- This Reduced data transfer → Minimizes bandwidth usage and speeds up analysis.

**but**

- Increased complexity → Requires sophisticated infrastructure and maintenance.
- Security restrictions → Limited user access and constrained software environments.
- Reduced flexibility → Users may not have full control over the computational setup.

**Some facilities already offer some computational environments to enable in-place data analysis.**

**This could be a viable solution but people may want to deal with more than one survey at the time, so the problem is still out there**
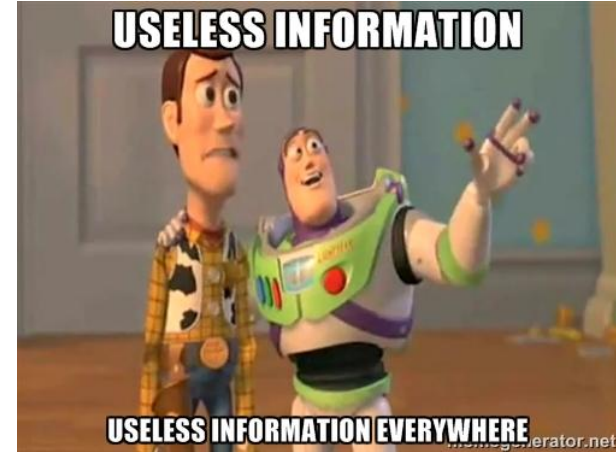
# Data Tsunami Mitigation

**Not all the users are interested to the same data**

Just to make the easiest example:

if you want to study galaxies you don't need to retrieve data of Stars

# I don't want the universe all at once

For a typical Machine Learning experiment (actually for any kind of science) what you need is of course a **subset of data points** and a **subset of the information** available for that point.

**I don't need (at least usually) to download everything is available in the archive.**

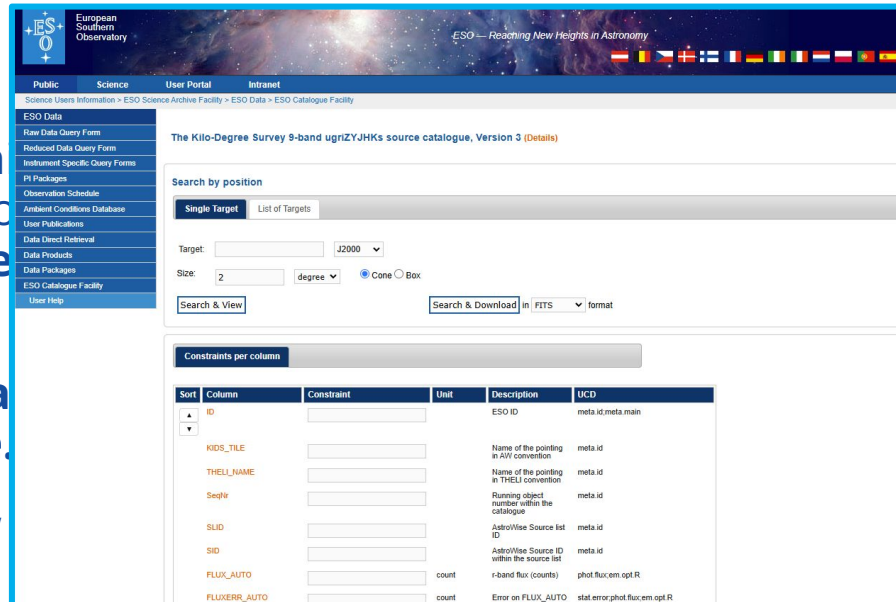**I need services that allow me to select only the data I need.**

# I don't want the universe all at once

For a typical Machine Learn[ing] (or any kind of science) what yo[u need is] **[a number] of data points** and a **subse[t]** [of the information] available for that point.

**I don't need (at least usua[lly] [all that] is available in the archive.**

**I need services that allow [me to take only what I] need.**



ESO portal, just as an example:
I have to download every single column and I have very limited possibility to perform a complex query
Even if a lot of survey are on the portal they do not "talk" one with the others

# I don't want the universe all at once

It should be obvious but:

**Data quality flags are really crucial when dealing with machine learning**

It is important that they are available and well documented

# I don't want the universe all at once

When comparing objects across multiple surveys, users need to cross-match them based on celestial coordinates.

This can happen on different levels:

- I have a catalogue of sources and I want the matching entries (**ideal solution would be to run the cross-match on the server rather than download everything and perform the match on my side**)
- I want sources with a given counterpart from a different survey, since this task could be often usual **add a pre calculated ID referencing to some of the most important surveys** would be useful sparing a lot of computation

# I don't want the universe all at once

## PHOTOMETRIC REDSHIFTS FOR QUASARS IN MULTI-BAND SURVEYS

M. BRESCIA[1,2], S. CAVUOTI[2], R. D'ABRUSCO[3], G. LONGO[2,4], AND A. MERCURIO[1]
[1] INAF-Astronomical Observatory of Capodimonte, via Moiariello 16, I-80131 Napoli, Italy; brescia@oacn.inaf.it
[2] Department of Physics, University Federico II, via Cinthia 6, I-80126 Napoli, Italy
[3] Harvard-Smithsonian Center for Astrophysics, 60 Garden Street, Cambridge, MA 02138, USA
Received 2013 February 28; accepted 2013 May 23; published 2013 July 17

…veys, users need …dinates.

**4 surveys involved**

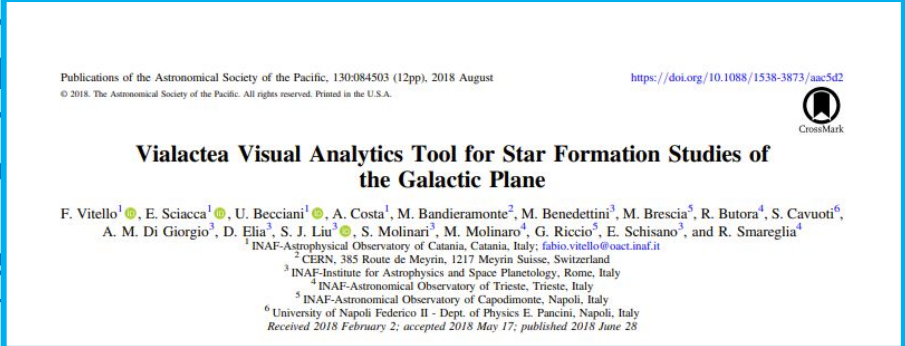- I have a catalogue of sources and I want the matching entries (**ideal solution would be to run th… cross-match on the server rather than d… everything and perform the match on m…**

- I want sources with a given counterpart fro… survey, since this task could be often usual… **calculated ID referencing to some of the… important surveys** would be useful sparin… computation

## Vialactea Visual Analytics Tool for Star Formation Studies of the Galactic Plane

F. Vitello[1], E. Sciacca[1], U. Becciani[1], A. Costa[1], M. Bandieramonte[2], M. Benedettini[3], M. Brescia[5], R. Butora[4], S. Cavuoti[6], A. M. Di Giorgio[3], D. Elia[3], S. J. Liu[3], S. Molinari[3], M. Molinaro[4], G. Riccio[5], E. Schisano[3], and R. Smareglia[4]
[1] INAF-Astrophysical Observatory of Catania, Catania, Italy; fabio.vitello@oact.inaf.it
[2] CERN, 385 Route de Meyrin, 1217 Meyrin Suisse, Switzerland
[3] INAF-Institute for Astrophysics and Space Planetology, Rome, Italy
[4] INAF-Astronomical Observatory of Trieste, Trieste, Italy
[5] INAF-Astronomical Observatory of Capodimonte, Napoli, Italy
[6] University of Napoli Federico II - Dept. of Physics E. Pancini, Napoli, Italy
Received 2018 February 2; accepted 2018 May 17; published 2018 June 28

**7 surveys involved**

# I don't want the universe all at once

In order to train a deep learning model **we need the images.**

**Not all surveys provide images and even in that case for most of them it is not straightforward to retrieve them**

We are often interested in small cutout surrounding the centroid of the object but we had to download entire plates and crop them locally, leading to significant inefficiencies. Users have access to raw, full-frame data, which can be useful for context and custom processing.

Bandwidth-heavy: Downloading entire plates is costly in terms of storage and data transfer. Although it would be computationally expensive for the server this effort in terms of computing power would be balanced from the limited usage of bandwidth.

Implementing a server-side cutout service that allows users to request only the relevant portion of an image directly from the archive. Some archives like SDSS and Pan-STARRS offer this, but it is not universal and they are implemented in a very different way and they could not be exactly what a user need (going back to the initial problem).

# I don't want the universe all at once

A&A 666, A171 (2022)
https://doi.org/10.1051/0004-6361/202243900
© L. Doorenbos et al. 2022

**Astronomy & Astrophysics**

## ULISSE: A tool for one-shot sky exploration and its application for detection of active galactic nuclei

Lars Doorenbos, Olena Torbaniuk, Stefano Cavuoti, Maurizio Paolillo, Giuseppe Longo, Massimo Brescia, Raphael Sznitman, and Pablo Márquez-Neila

**~100k 73x73px 3-band thumbnails**

I REPEAT, OPEN UP AND RELEASE THE DATA!

DATABASE

A&A, 687, A246 (2024)
https://doi.org/10.1051/0004-6361/202450166
© The Authors 2024

**Astronomy & Astrophysics**

## Identification of problematic epochs in astronomical time series through transfer learning

Stefano Cavuoti, Demetra De Cicco, Lars Doorenbos, Massimo Brescia, Olena Torbaniuk, Giuseppe Longo, and Maurizio Paolillo

**~1M 51x51px single band thumbnails**

THE ASTROPHYSICAL JOURNAL, 977:131 (22pp), 2024 December 10
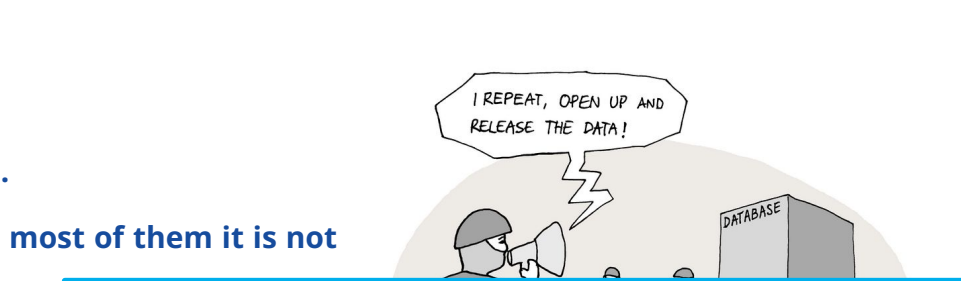© 2024. The Author(s). Published by the American Astronomical Society.
OPEN ACCESS
https://doi.org/10.3847/1538-4357/ad8bbe

## Galaxy Spectroscopy without Spectra: Galaxy Properties from Photometric Images with Conditional Diffusion Models

Lars Doorenbos, Eva Sextl, Kevin Heng, Stefano Cavuoti, Massimo Brescia, Olena Torbaniuk, Giuseppe Longo, Raphael Sznitman, and Pablo Márquez-Neila
AIMI, ARTORG Center, University of Bern, Murtenstr. 50, CH-3008 Bern, Switzerland; lars.doorenbos@unibe.ch
Universitäts-Sternwarte, Fakultät für Physik, Ludwig-Maximilians Universität München, Scheinerstr. 1, 81679 München, Germany; sextl@usm.lmu.de
INAF—Astronomical Observatory of Capodimonte, Salita Moiariello 16, I-80131 Napoli, Italy
INFN—Sezione di Napoli, via Cinthia 9, 80126 Napoli, Italy
Department of Physics, University Federico II, Strada Vicinale Cupa Cintia, 21, 80126 Napoli, Italy
Department of Physics and Astronomy "Augusto Righi," University of Bologna, via Piero Gobetti 93/2, 40129 Bologna, Italy
Received 2024 June 24; revised 2024 October 16; accepted 2024 October 25; published 2024 December 9

**~300k 64x64px 5-band thumbnails**

# I Don't Want to Miss a Thing

On the other hand I need to have control on missing data. Usually in astronomical DB a value could be missing for mainly three reason:
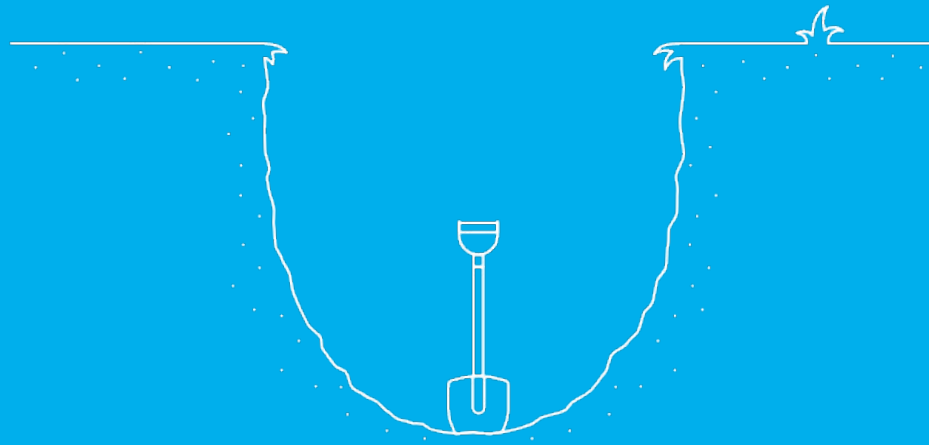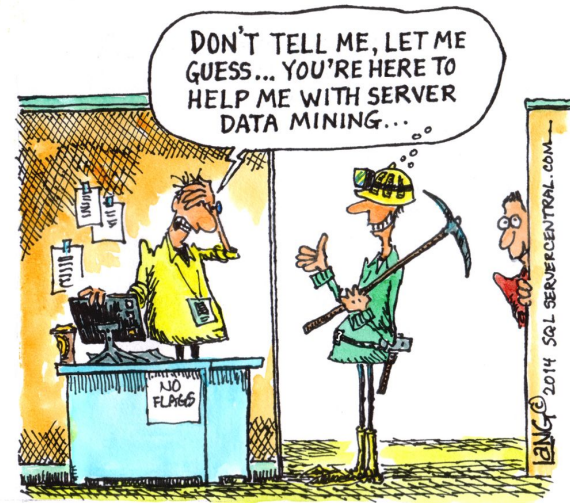
A. Observed but not detected
B. Observed but the object has been masked for a problem
C. Not observed

It is really important to distinguish at least between A. & B. from C. since I would deal with them in a very different way.

# Thank you for your attention!

# Possible Discussion Points

- computational environments to enable in-place data analysis
- complex query infrastructure
- importance of data quality flags
- in-place cross-match
- precomputed cross-match with surveys
- on-the-fly cutout service
- missing data information