# Applications of multiple DBMSs and algorithms for time-domain astronomy

Min-Su Shin (Korea Astronomy and Space Science Institute) msshin@kasi.re.kr
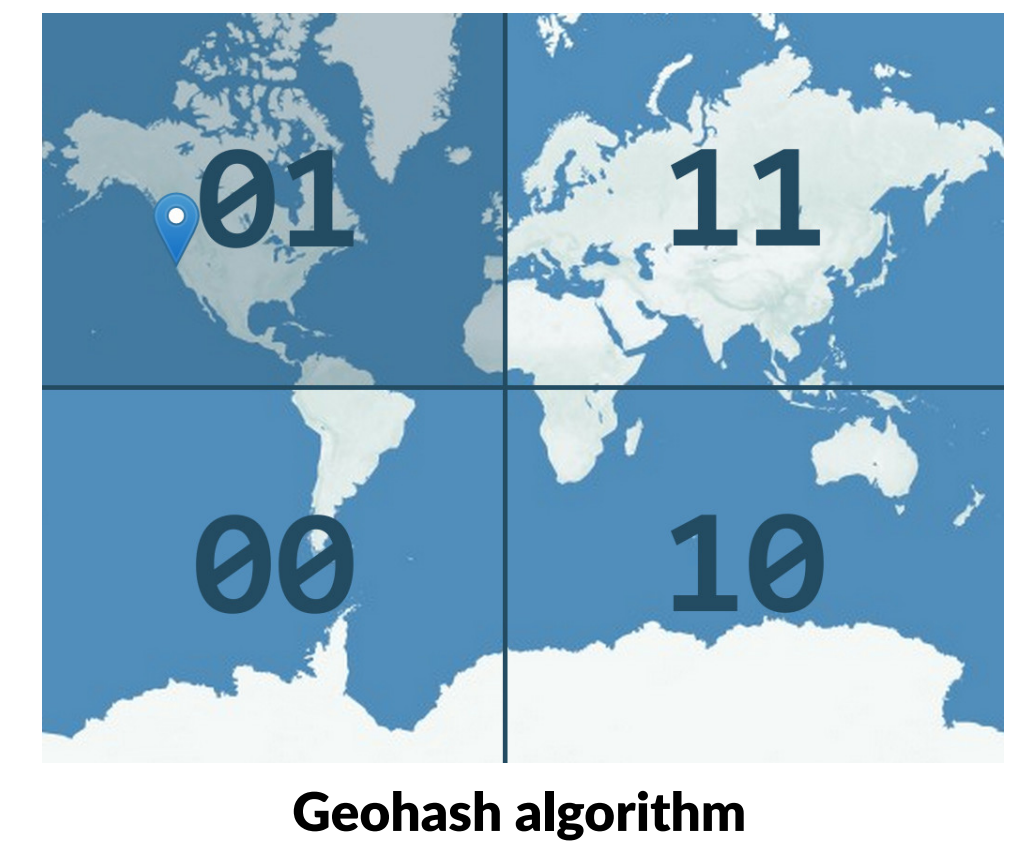
KASI — Korea Astronomy & Space Science Institute

We use multiple DMBSs and algorithms in the follow-up target selection step, processing follow-up observation data, and cataloging reduced data and light curves. The current system is used in our pilot program of time-domain follow-up observations.

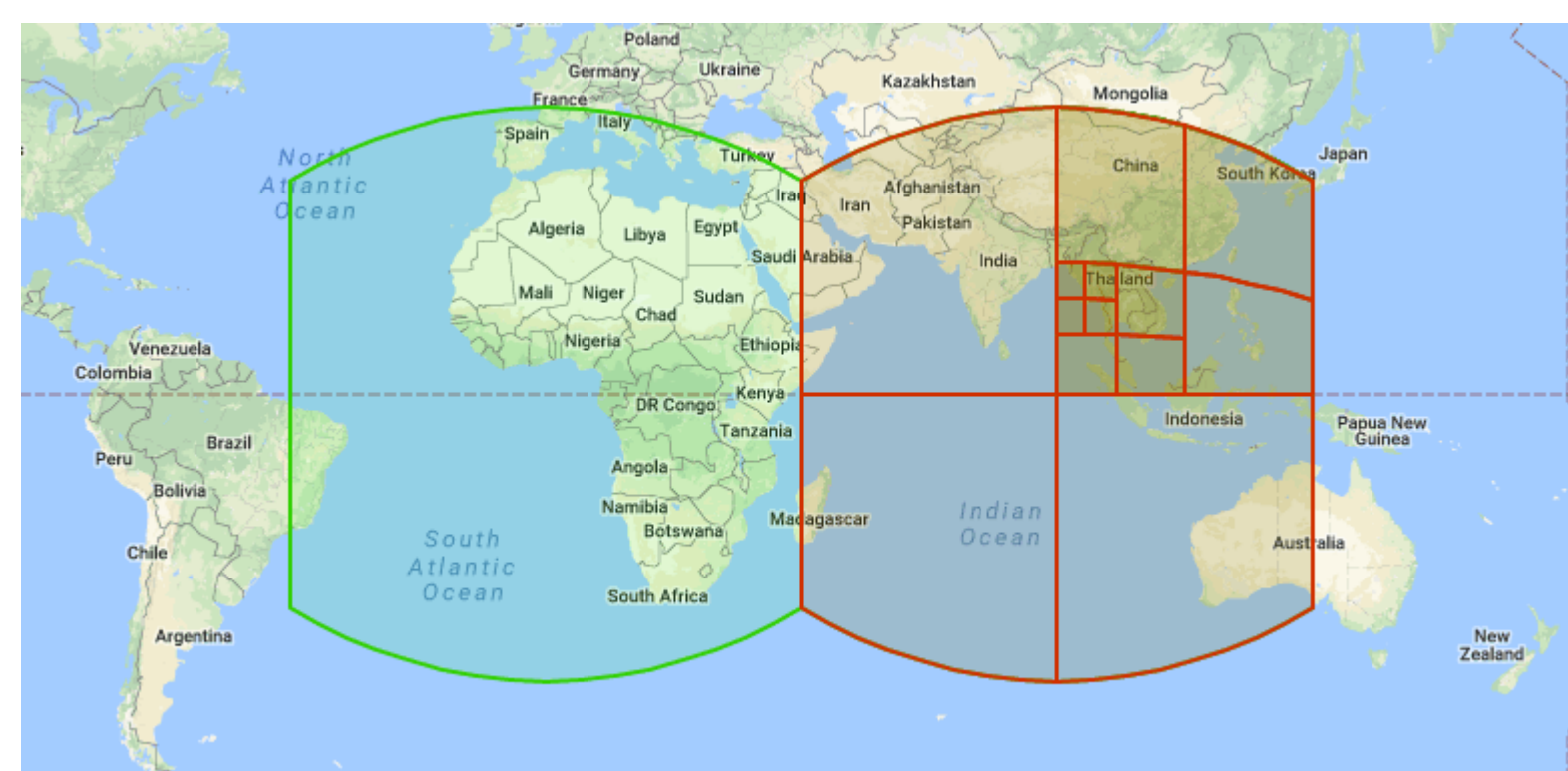## 1. Redis: GeoSet (a sorted set with latitude and longitude)

**Low-latency in-memory spatial data store** for astronomical coordinates.

−Modified version of the Redis to store custom catalogs with coordinates for follow-up target selection or local catalog search purposes.
−Typical search response time ~ microseconds to milliseconds thanks to the geohash algorithm.

Local custom catalog database
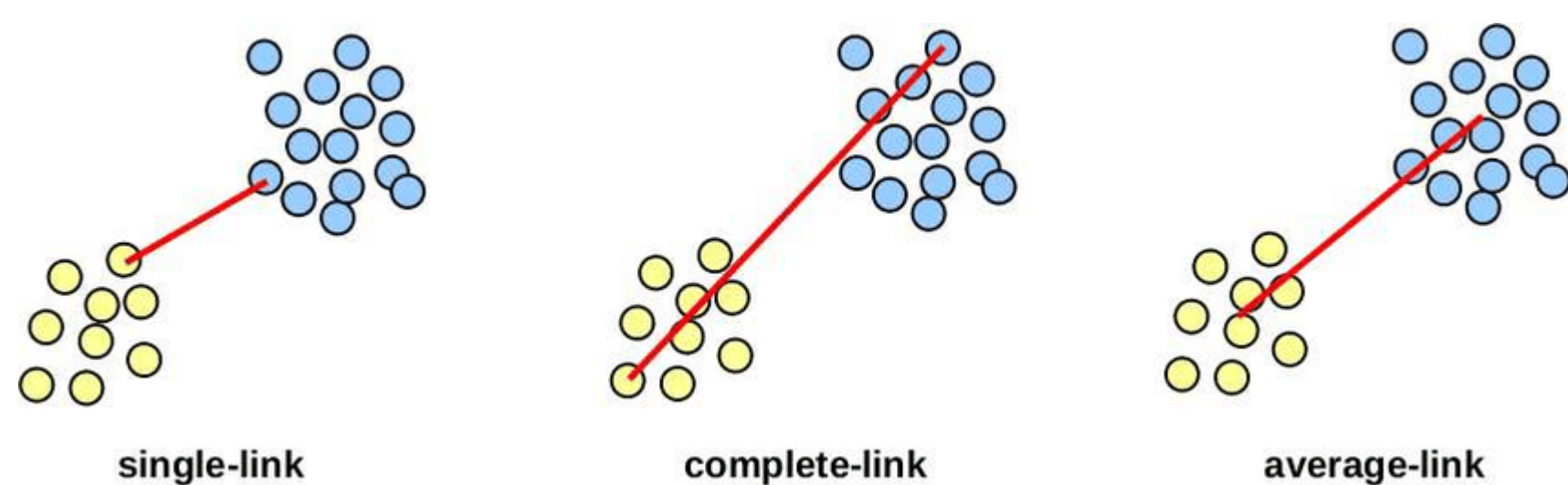
redis

Geohash algorithm

## 2. Google's S2 Geometry library

Google's S2 Geometry library

single-link    complete-link    average-link

Constructing light curves and physical object catalogs by single-linkage (i.e., friends-of-friends) clustering of detected sources with helps of **Google's S2 Geometry library for fast spatial search**.

− Input: catalogs of detected sources with their positions in optical observations for given observation fields.
− Output: light curves and physical objects defined as grouped sources with a given linking angular distance in the single-linkage single-level clustering.
− Our program written in C++ using Google's S2 Geometry library for indexing detected sources and searching the nearest neighbor in the single-linkage clustering.

## 3. ClickHouse: column-oriented DBMS for source and object catalogs

Our requirements for data store of source, object, and image catalogs:
− Horizontally scalable (i.e. sharding) with commodity hardware.
− SQL-like query support.
− Reasonable data ingestion performance.
− Fast search query performance with group by observation field names or for spatiotemporal constraints.

ClickHouse    RethinkDB    Vitess

Our consideration and test of three open-source systems: **ClickHouse, RethinkDB, and Vitess**.
− Vitess: (pros) the most powerful choice with the broad supports of MySQL features and sharding with various custom rules, (cons) requirements of the k8s cluster and large resources.
− RethinkDB: (pros) automatic sharding and fast response in ingestion and search, (cons) weak community development and large resources.
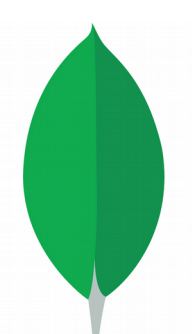− ClickHouse: (pros) the efficient usage of resources, supports of sharding with SQL-like languages, (cons) the limited support of sharding rules and no geodata features. → We adopt the **ClickHouse as the main store of source, object, and image catalogs for their fixed schema and infrequent usage of entire columns**.

## 4. MongoDB: document-oriented DBMS for light curves

mongoDB

− Why?
a. different sizes and contents of light curves well matched to documents stored in MongoDB.
b. supports of sharding and geodata in document-format data.

## 5. Plan

−Continuous tests with the current data reduction process and stores for our pilot observation program until Sep. 2020.
− Expanding the current configuration of DBMS nodes in a public cloud from 4 nodes to 8 nodes for ClickHouse and MongoDB, respectively.
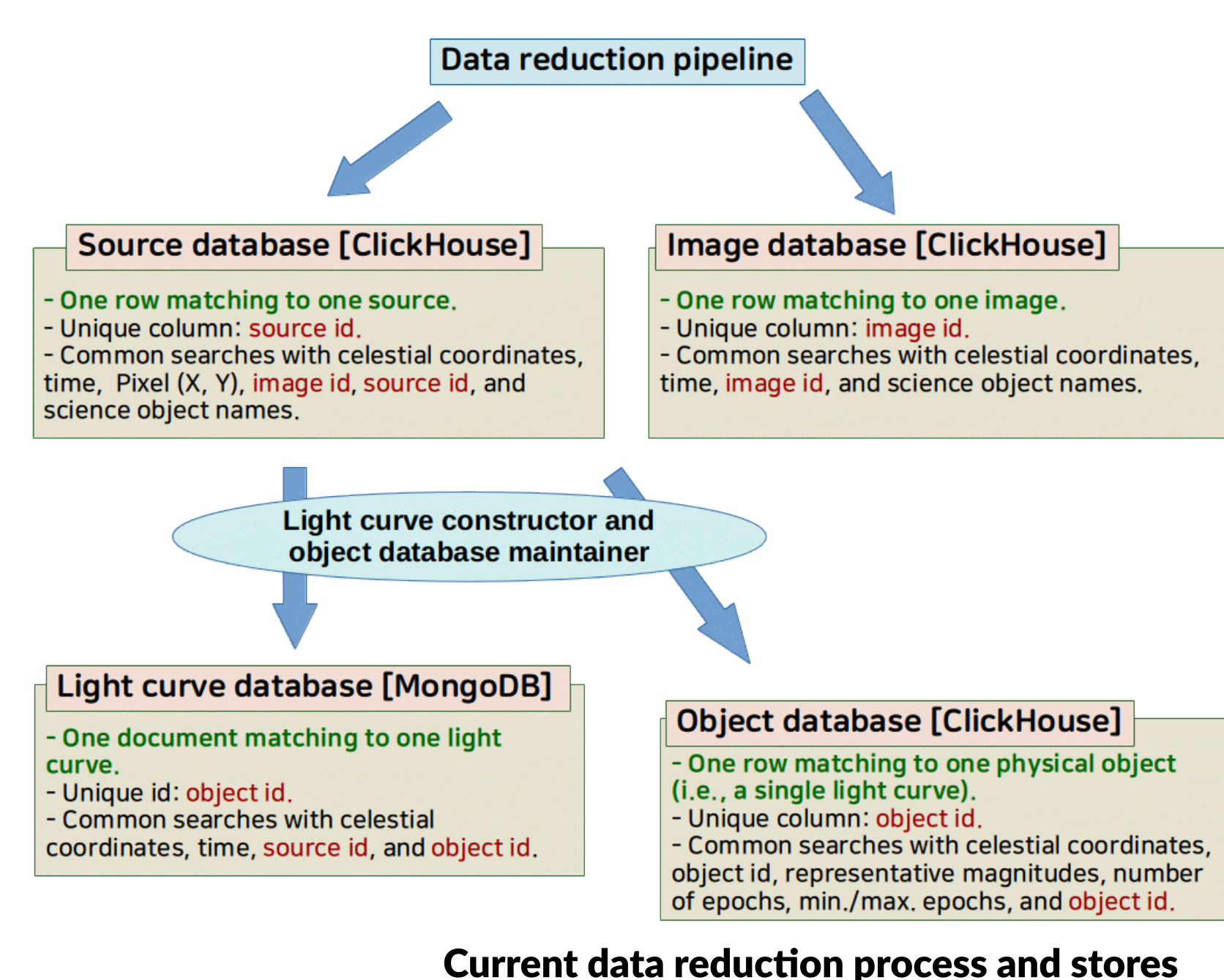
```
source_id CHAR(40) NOT NULL,
seq INT UNSIGNED NOT NULL,
# source sequence number (for a specific amp)
filter CHAR(1),
image_id CHAR(40) NOT NULL,
sci_obj_name CHAR(128),
mjd DOUBLE NOT NULL,
# MJD (day)
x DOUBLE NOT NULL,
# X image (SEXTRACTOR KEYWORDS)
y DOUBLE NOT NULL,
# Y image (SEXTRACTOR KEYWORDS)
ra_deg DOUBLE NOT NULL,
# RA world (SEXTRACTOR KEYWORDS; degree)
dec_deg DOUBLE NOT NULL,
# DEC world (SEXTRACTOR KEYWORDS; degree)
ra_dec_point POINT NOT NULL SRID 4326,
mag_auto DOUBLE,
# calibrated Mag. auto (SEXTRACTOR KEYWORDS; mag_auto)
magerr_auto DOUBLE,
# uncertainty of mag_auto (SEXTRACTOR KEYWORDS; mag_auto uncertainty)
bkg DOUBLE,
# background at centroid position (SEXTRACTOR KEYWORDS; ADU)
fwhm DOUBLE,
# FWHM assuming a gaussian core (SEXTRACTOR KEYWORDS; pixel)
ellipticity DOUBLE,
# ELLIPTICITY (SEXTRACTOR KEYWORDS)
class_star DOUBLE,
# S/G classifier output (SEXTRACTOR KEYWORDS)
sex_flag SMALLINT UNSIGNED,
# SExtractor extraction flags (SEXTRACTOR KEYWORDS)
mag_map DOUBLE,
# photometrically calibrated best mag.
magerr_map DOUBLE,
# (derived and corrected by MAP; magnitude)
magerr_map DOUBLE,
# uncertainty of mag_map (magnitude)
ap_map DOUBLE,
# aperture diameter for mag_map (pixel)
refmag_map DOUBLE,
# magnitude with Max. AP via MAP (magnitude)
avg_delta_m DOUBLE,
# average of photometric calibration delta_m
std_delta_m DOUBLE,
# standard deviation of photometric calibration delta_m
skew_delta_m DOUBLE,
# skewness of photometric calibration delta_m
source_reliability DOUBLE,
photometry_reliability_n INT
photometry_reliability_1 DOUBLE,
photometry_reliability_2 DOUBLE,
photometry_reliability DOUBLE
```
**Source catalog table**

Data reduction pipeline

**Source database [ClickHouse]**
- One row matching to one source.
- Unique column: source id.
- Common searches with celestial coordinates, time, Pixel (X, Y), image id, source id, and science object names.

**Image database [ClickHouse]**
- One row matching to one image.
- Unique column: image id.
- Common searches with celestial coordinates, time, image id, and science object names.

**Light curve constructor and object database maintainer**

**Light curve database [MongoDB]**
- One document matching to one light curve.
- Unique id: object id.
- Common searches with celestial coordinates, time, source id, and object id.

**Object database [ClickHouse]**
- One row matching to one physical object (i.e., a single light curve).
- Unique column: object id.
- Common searches with celestial coordinates, object id, representative magnitudes, number of epochs, min./max. epochs, and object id.

**Current data reduction process and stores**

{
_id: ObjectID(),
(given by mongoDB)
object_id: KMTNJ083025.53-070822.5,
(produced by the constructor as the mean of RA and DEC)
ra: RA,
(mean value for sources with source_reliability >= 19;
degree)
dec: DEC,
(mean value for sources with source_reliability >= 19;
degree)
num_obs_lc: 111,
(number of observation epochs given to the linking
procedure)
num_obs: 107,
(number of observation epochs with reliable sources, i.e.,
source_reliability >= 19, used for the linking procedure)
obs: [
{source_id: 20190123032334p54_CTIO_09_34p73_32p34,
filter: V, mjd: 653109.345572, ra: 34.567456, dec: -
74.341096,
mag_map: 16.564, magerr_map: 0.014,
mag_auto: 16.780, magerr_auto: 0.023,
sex_flag: 0,
source_reliability: 65,
photometry_reliability_n: 23,
photometry_reliability_1: 0.003,
photometry_reliability_2: 0.032,
photometry_reliability: 0.931},
{source_id: 20190123032423p26_CTIO_09_34p71_32p32,
filter: V, mjd: 653109.985572, ra: 34.657456, dec: -
74.341086,
mag_map: 16.545, magerr_map: 0.0143,
mag_auto: 16.791, magerr_auto: 0.021,
sex_flag: 0,
source_reliability: 86,
photometry_reliability_n: 36,
photometry_reliability_1: 0.0013,
photometry_reliability_2: 0.0209,
photometry_reliability: 0.71}

process_datetime: date and time in UTC
}
**Light curve example**