



NOVEL MATERIALS DISCOVERY

BBDC FHI Meeting

NOMAD: Data & Pipelines

14.02.2017

Fawzi Mohamed FHI Berlin



NOVEL MATERIALS DISCOVERY

Data & Pipelines

- Raw Data
 - Zip archives (Bag it format)
 - Text or binary data in various formats
- Parsed Data
 - Json or HDF5
 - Metadata representation
- Normalized Data
 - (Json) HDF5 (+ Parquet?)
 - Metadata representation
- Parsing
 - Raw Data → Parsed Data
- Normalization
 - Parsed Data → Normalized Data
- Query
 - Normalized Data → Query result



NOVEL MATERIALS DISCOVERY

Metadata

- ~dictionary, ontology, schema
- Unique name for each entry
- Sections
 - Collect values that belong together (Dictionary)
 - Nested (form a tree)
- Concrete values
 - Multidimensional array of
 - Boolean, string, int/long, real/float
 - Or reference to a section (relational model)



NOVEL MATERIALS DISCOVERY

MetaData

```
section_run
  program_name      FHI-aims
  program_version   081912
section_system
  simulation_cell  [[1.4e-9 ...]]
  atom_positions    [[0.0, ...] ...]
  atom_labels       ["Cu", ...]
section_method
  basis_set        fhi_aims_tight
  XC_method        DFT_GGA_PBE
section_single_configuration_calculation
  section_scf_iteration
    energy_total_scf_iteration -1.326e-20
  section_scf_iteration
    energy_total_scf_iteration -1.344e-20
energy_total      -1.344e-20
```

Values: Data
Structure
and names: Metadata

SI Units:

- lengths: m
- energies: J
- ...

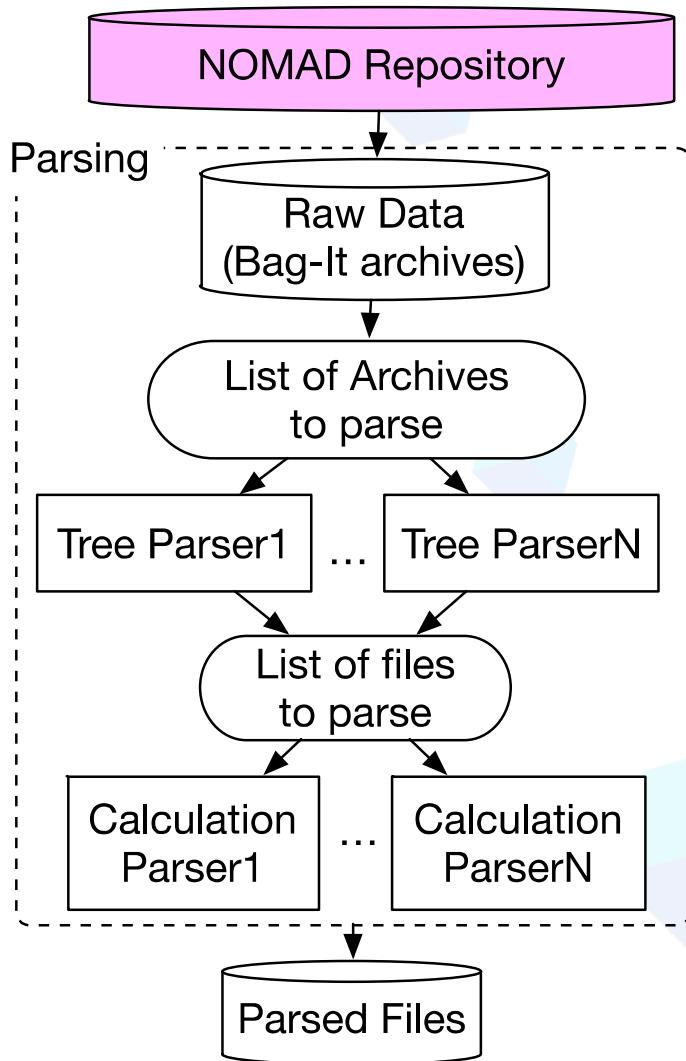


NOVEL MATERIALS DISCOVERY

Parse Events

- Open Section
- Close Section
- Add Value
- Partially ordered
- Fully ordered

Parsing



- Parsers **interpret** all calculation data
- **Organize** it according to the metadata structure
- Data not extracted is invisible
- Writing a parser cannot be automatized and requires a person with scientific knowledge
- **Parallel** execution
 - *Tree Parser* identifies the files
 - *Calculation Parser* performs the parsing and generates the parsed files
- Parsing is **pure**: the same version on the same data should give the same result

Parser Developers



Mikkel Strange



Micael Oliveira



Martina Stella



Massimo Riello



Wael Chibiani



Franz Knuth



Aliaksei Mazheika



Adam Fekete



Andrea Droghetti



Honghui Shang



Sebastian Alarcón Villaseca



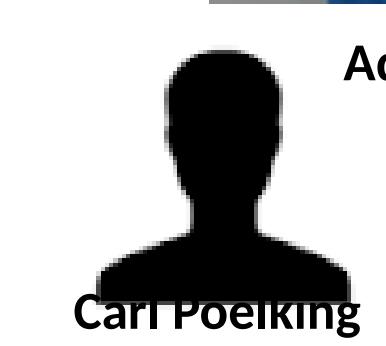
Daria Tomecka



Lorenzo Pardini



Fawzi Mohamed



Carl Poelking



Adriel Dominiguez



Rosenndo Valero



Ask Hjorth Larsen



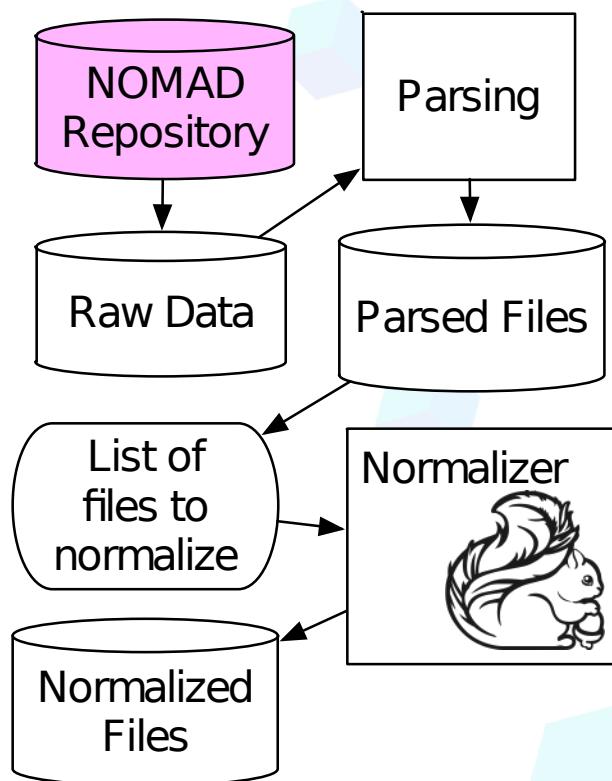
Lauri Himannen



Sami K. Kivistö



Henning Glawe



Standardization

- Parsers standardize the data,
- avoid the loss of information

Normalization

- The work packages define **derived quantities** (normalized representations,...)
- These can be generally useful for analysis or visualization
- Normalization is an infrastructure to apply automatically some transformations and store their result along with the parsed data



NOVEL MATERIALS DISCOVERY

NOMAD Archive in Numbers

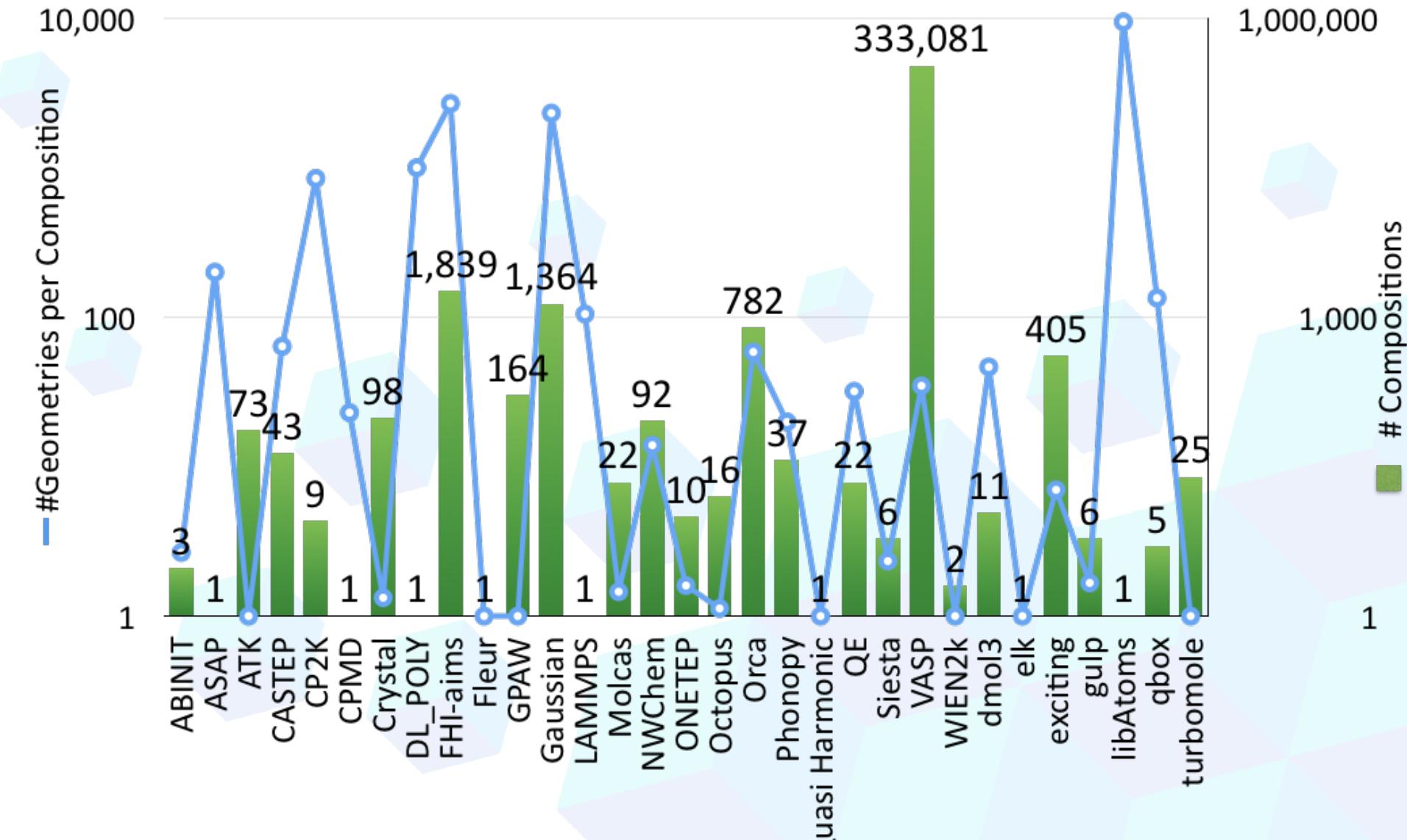
generated parsing
corresponding to
Normally the content of the zip archives are repetitive text files
with excellent
This produces
with data classified using

3761 Zip Archives
8 TeraBytes of compressed data.
compression ratio **60-80%**.
2.7 TeraBytes of parsed Hdf5 files
279 public metadata
1736 code specific metadata



NOMAD

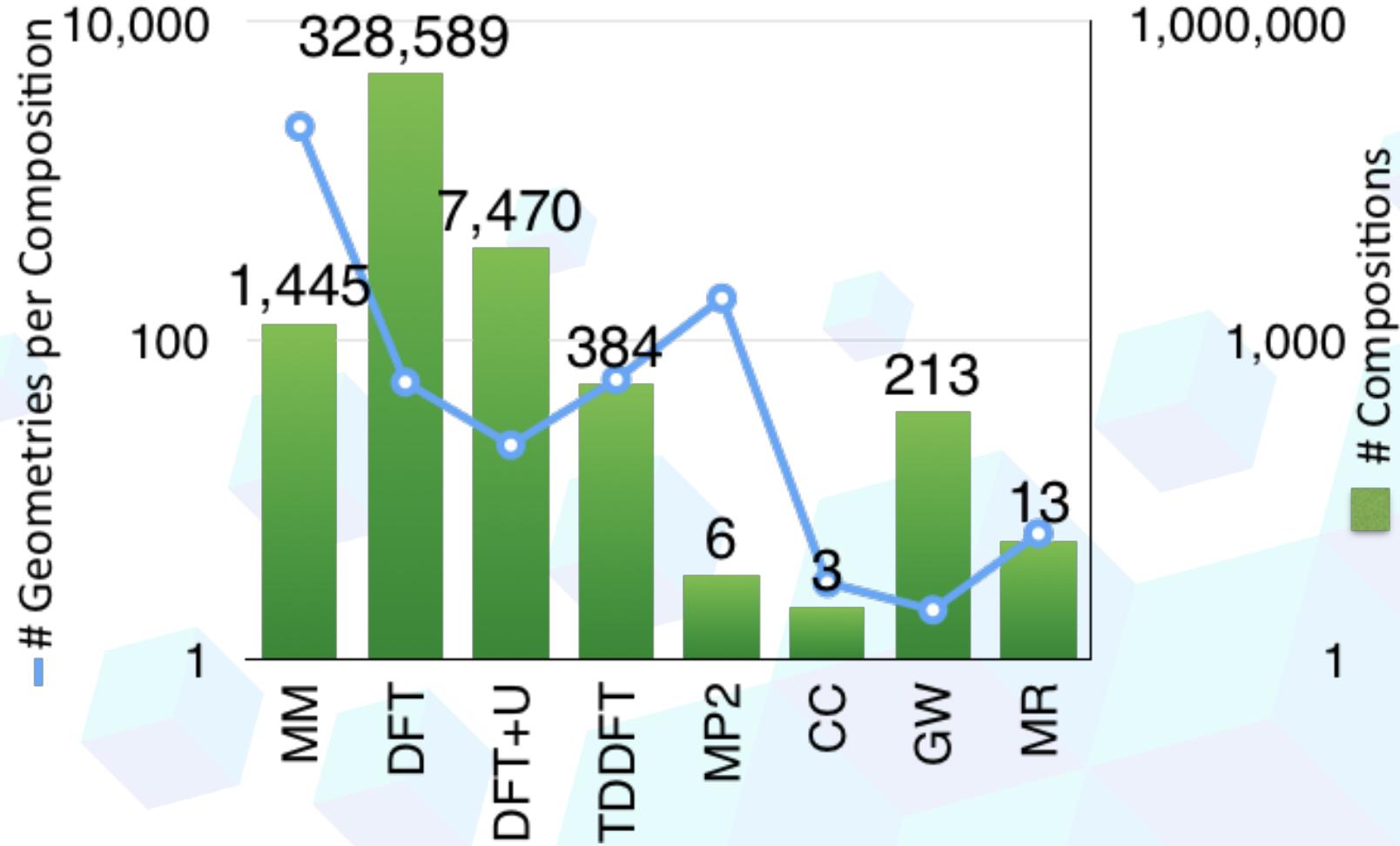
NOVEL MATERIALS DISCOVERY





NOMAD

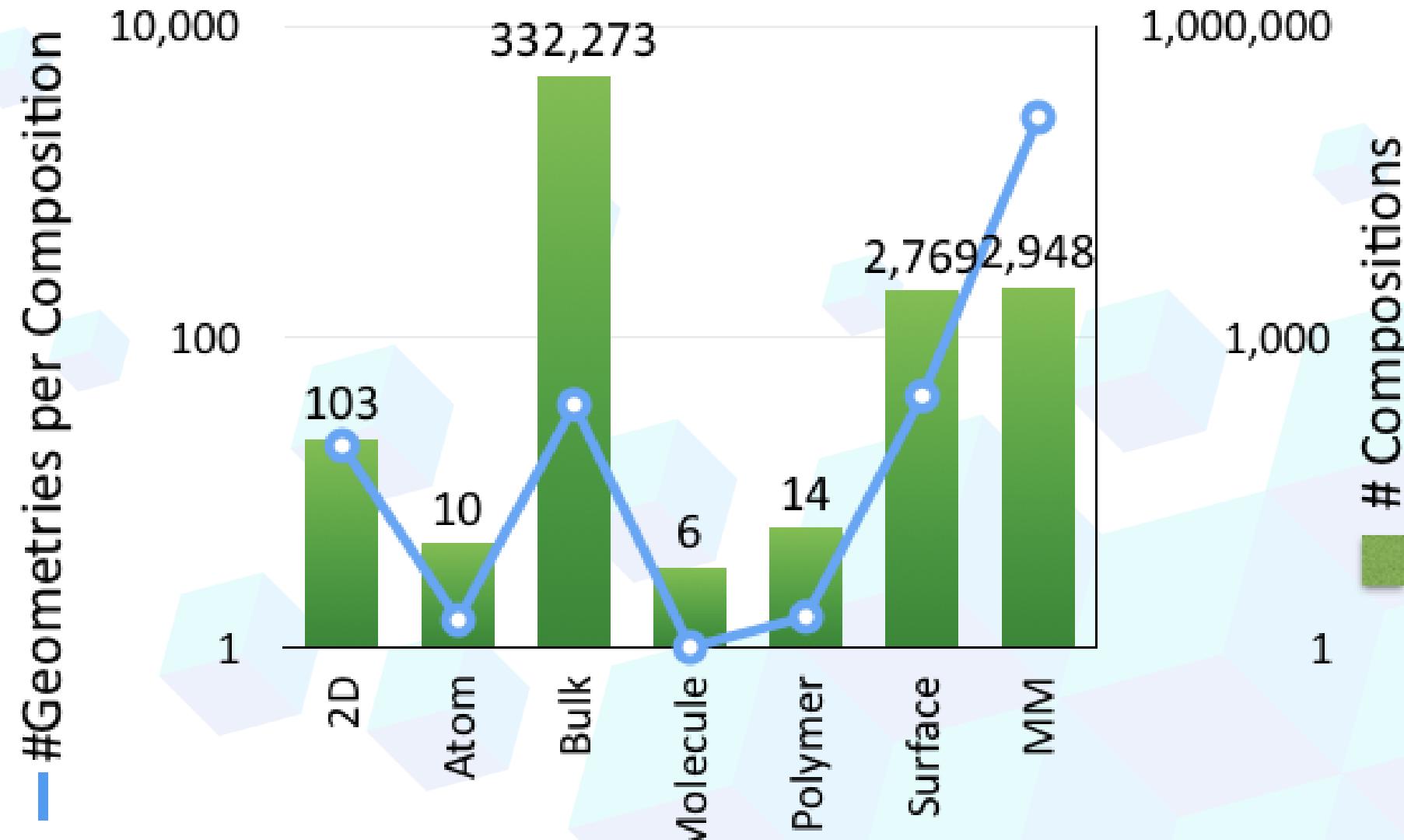
NOVEL MATERIALS DISCOVERY





NOMAD

NOVEL MATERIALS DISCOVERY





NOVEL MATERIALS DISCOVERY

Next Steps

- HDF5 python interface
- Support when analysis is preformed (Bug fixing)
- 10 extra Parsers April 2018 (M30)
- James Kermode (TINKER, NAMD, AMBER, Gromos, Gromacs, CHARMM)
- Aiida
- LM Suite (TB-LMTO-ASA) (KCL?),
- WEST



NOVEL MATERIALS DISCOVERY

NOMAD Metadata

Metadata:

ESCDF
HJSON
NIST, EMMC

OPTIMADE REST API

Normalization

279 public metadata

1736 code specific metadata



NOVEL MATERIALS DISCOVERY

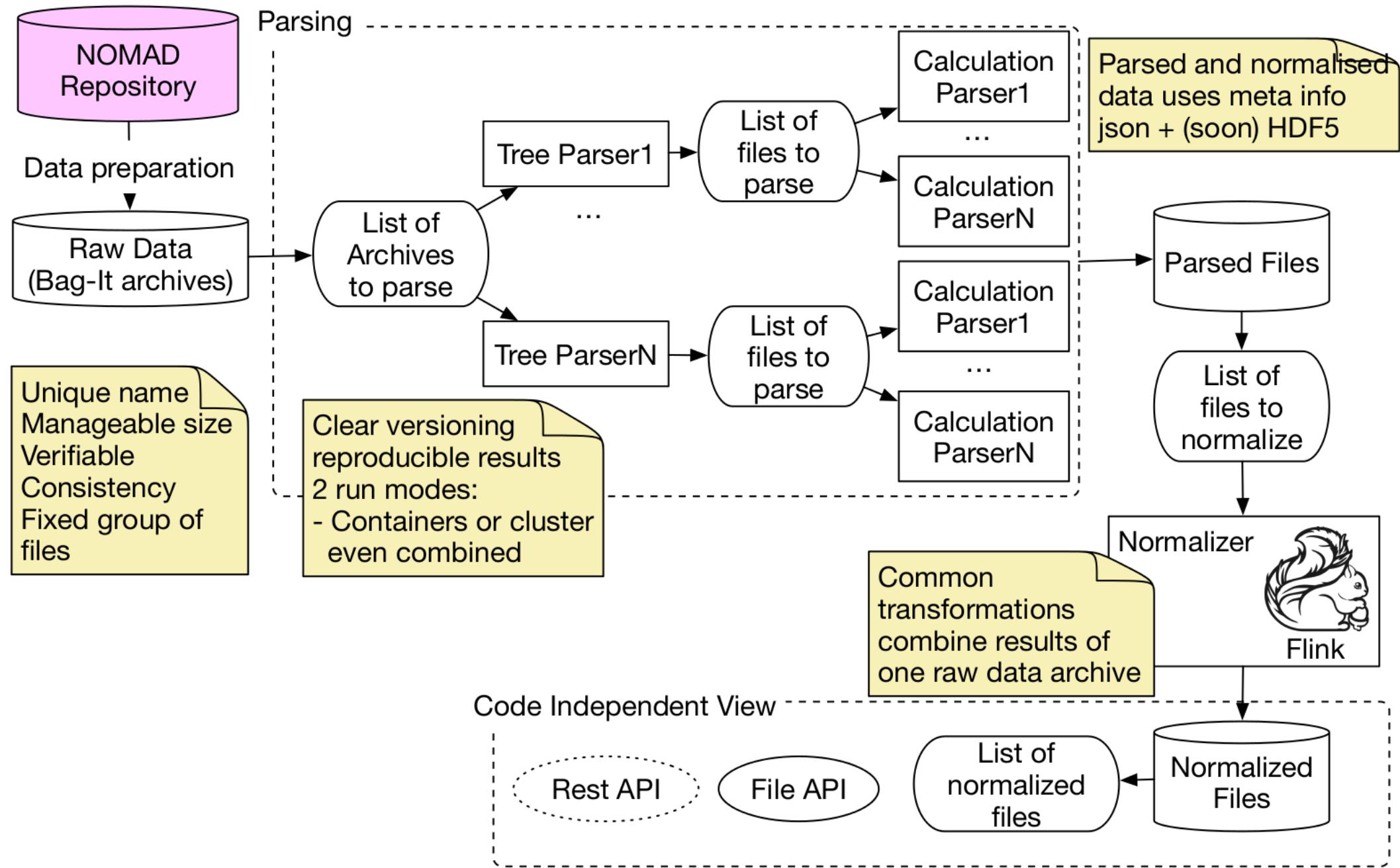
Normalization

- Messy part, might have dependencies

Example:

- Checksum for geometry
- Connection to repository
- Statistics

- Loop on archives/calculations/single point/goemetries
 - Apply transformation
 - Store result
- OpenContext/closeContext event





NOVEL MATERIALS DISCOVERY

Query

- Read only
- Loop on data, select subset
- Return nomad uri + data table