



NOVEL MATERIALS DISCOVERY

# WP1 Report

## The NOMAD Laboratory

EUSpec Meeting ESCDF further steps

Budapest 18 Oct. 2016

Fawzi Mohamed FHI Berlin

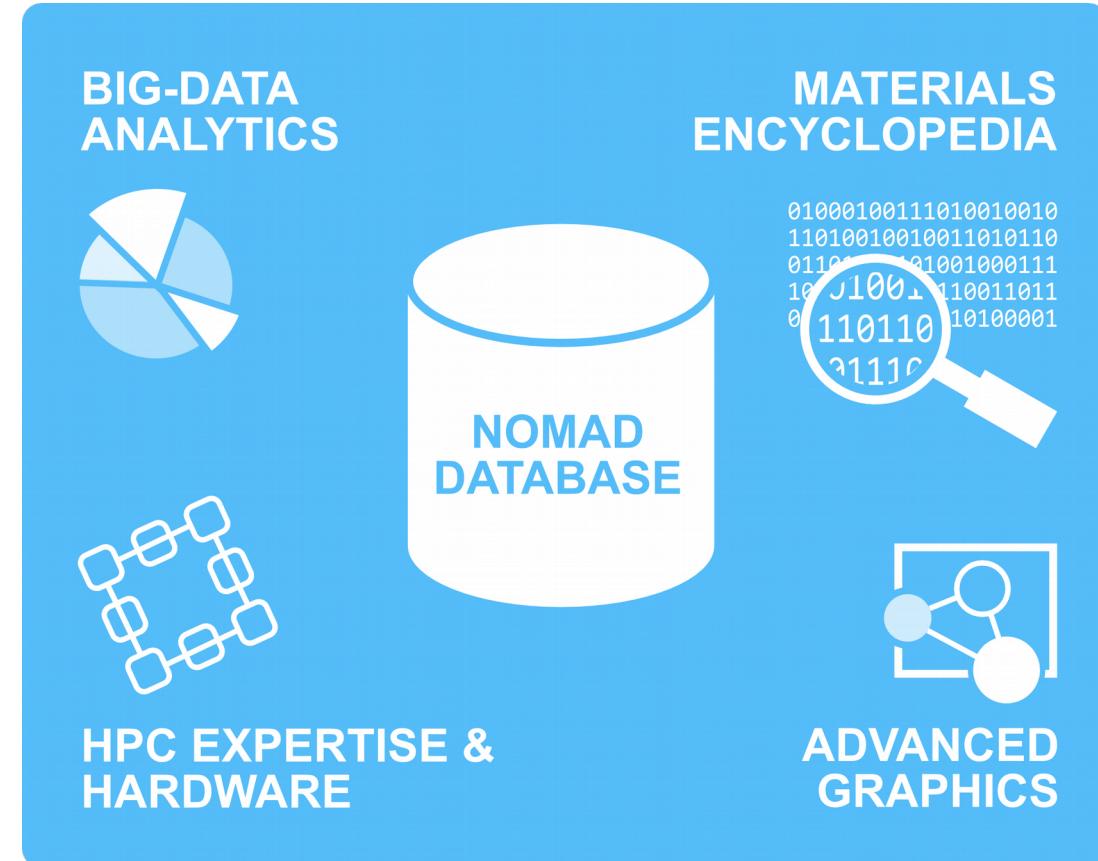


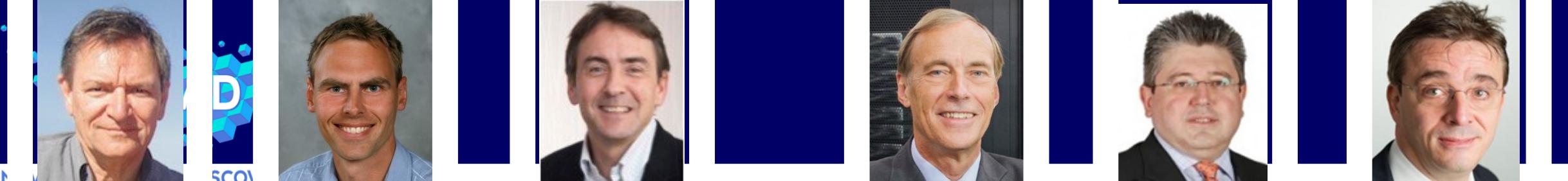
NOVEL MATERIALS DISCOVERY

## The Novel Materials Discovery (NOMAD)

Laboratory develops a Materials Encyclopedia and Big-Data Analytics and Advanced Graphics Tools for materials science and engineering.

Eight complementary computational materials science groups and four high-performance computing centers form the core of this Centre of Excellence.





**Matthias Scheffler**, FHI  
MPS, Berlin



**Angel Rubio** MPI  
MPSD, Hamburg



**Risto Nieminen**  
Aalto U. Helsinki



**Kristian Thygesen**  
Tech. U., Lyngby



**Ciaran Clissman**  
pintail Ltd.  
Dublin



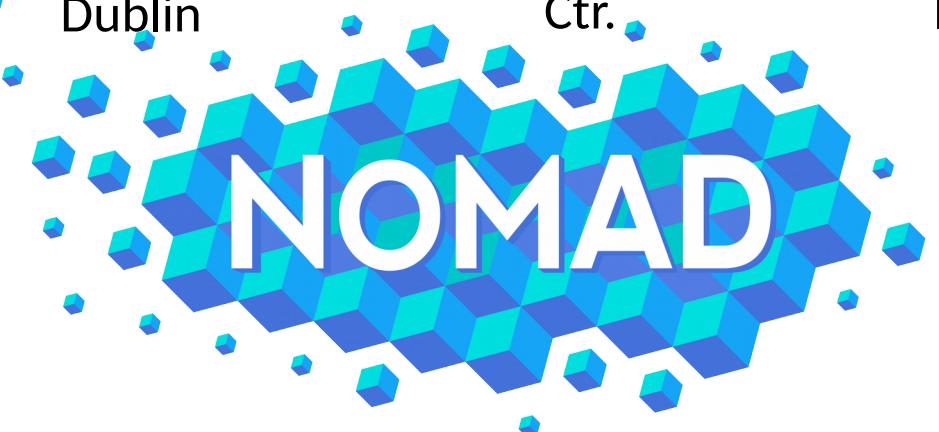
**Arndt Bode**  
Leibniz Comp.  
Ctr.



**Jose Maria Cela**, BSC,  
Barcelona



**Alessandro De Vita King's  
Col.** London



**NOVEL MATERIALS DISCOVERY**



**Kimmo Koski** CSC –  
IT Center Helsinki



**Francesc Illas** U.  
of Barcelona



**Stefan Heinzel**  
MPS Comp. & Data, Garching



**Daan Frenkel** U.  
Cambridge



NOVEL MATERIALS DISCOVERY

# NOMAD Archive: Goal





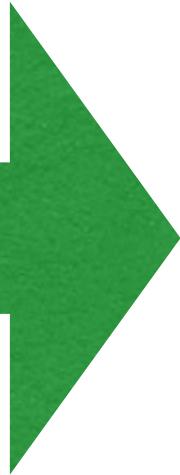
NOVEL MATERIALS DISCOVERY

NOMAD  
REPOSITORY

# NOMAD Archive: Goal

NOMAD Archive

Make data of the simulation codes  
used in the community accessible



WP2: NOMAD  
ENCYCLOPEDIA

WP3: NOMAD  
VISUALIZATION

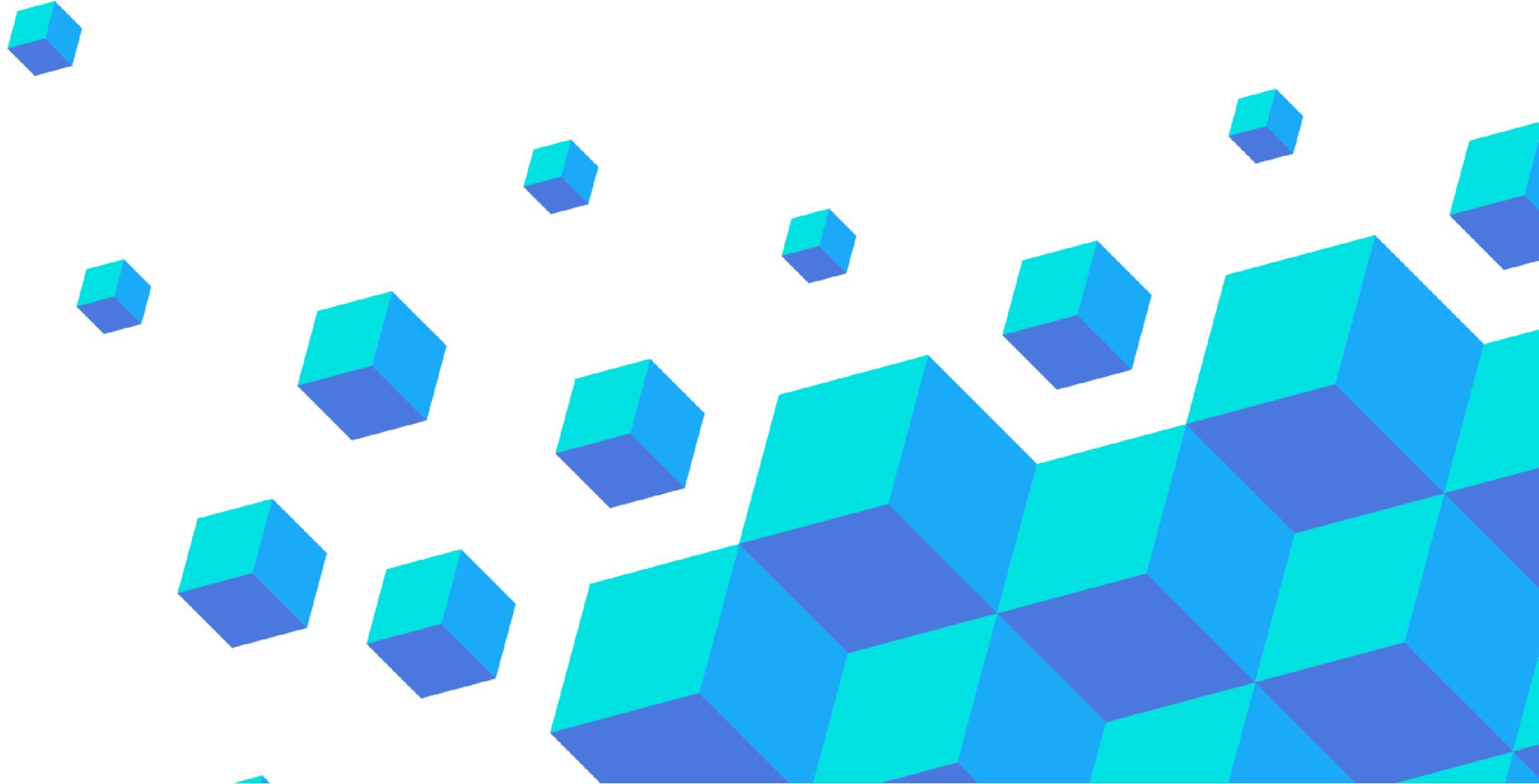
WP4: BIG DATA  
ANALYSIS

*External users*

NOMAD

NOVEL MATERIALS DISCOVERY

# NOMAD Archive: What is it

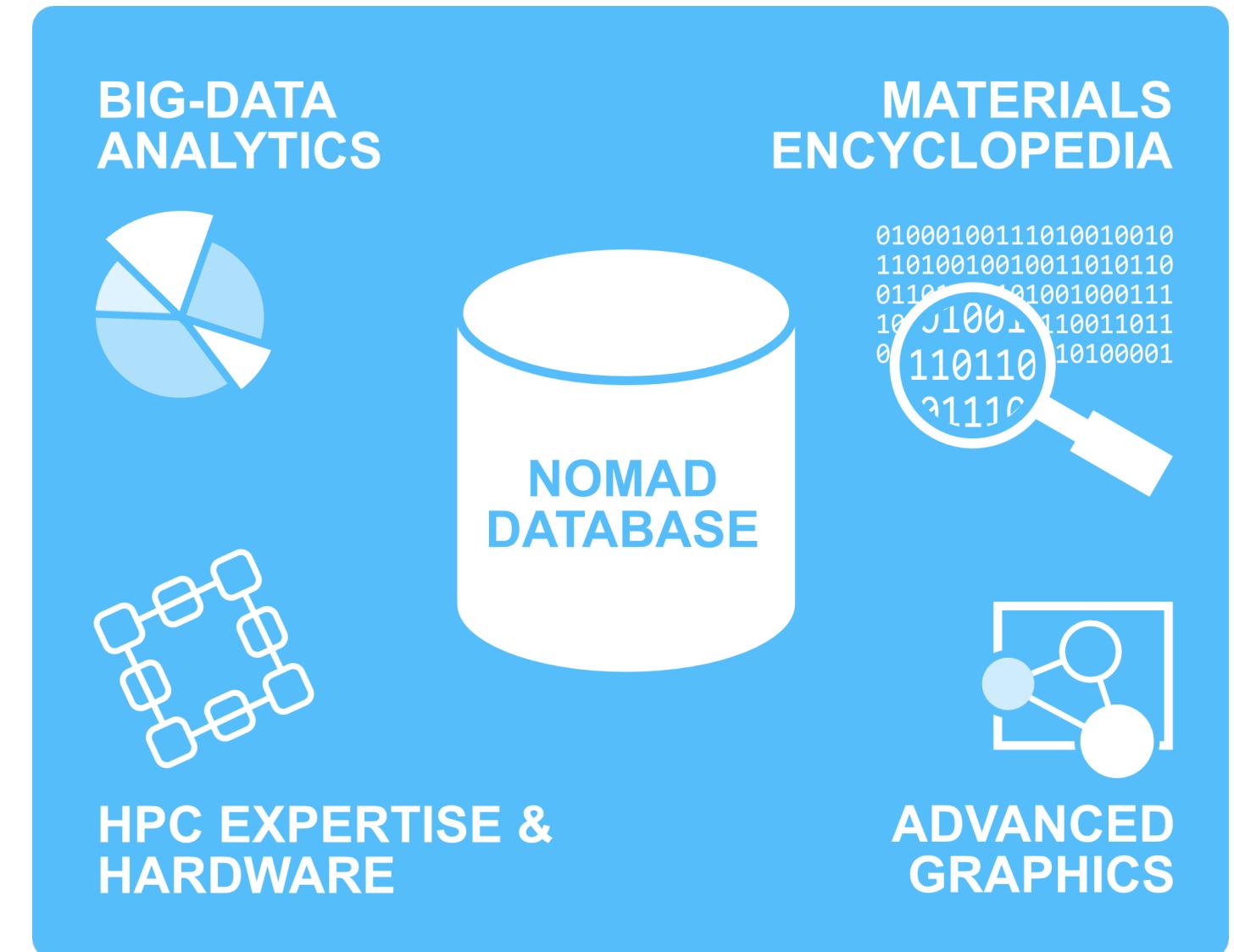




NOVEL MATERIALS DISCOVERY

The core of the  
NOMAD Database

# NOMAD Archive: What is it





# NOMAD Archive: What is it

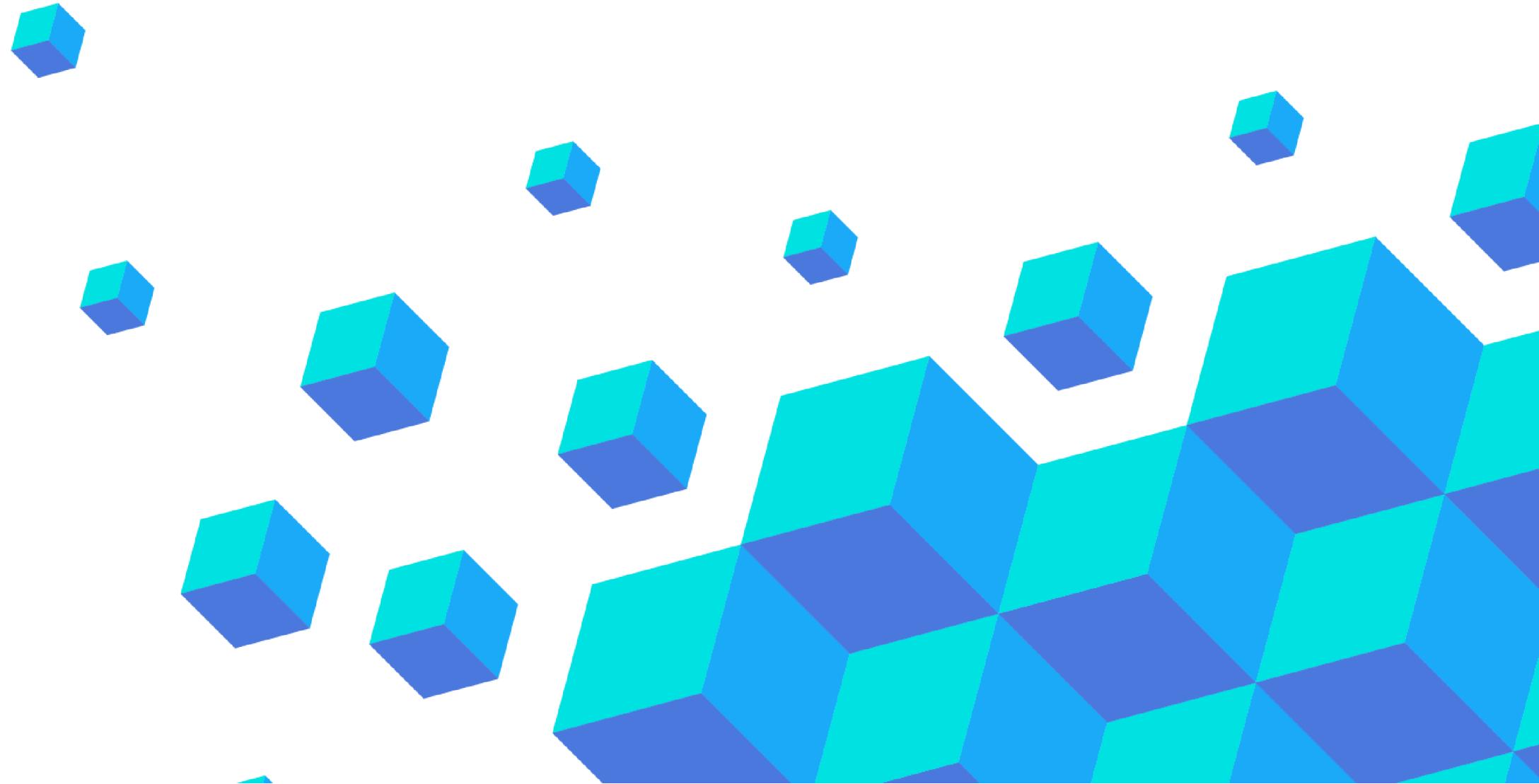
A growing collection of files

- Containing calculation data organized with **metadata**
- With clear **identifiers** to enable workflows and automatic processing
- Open access



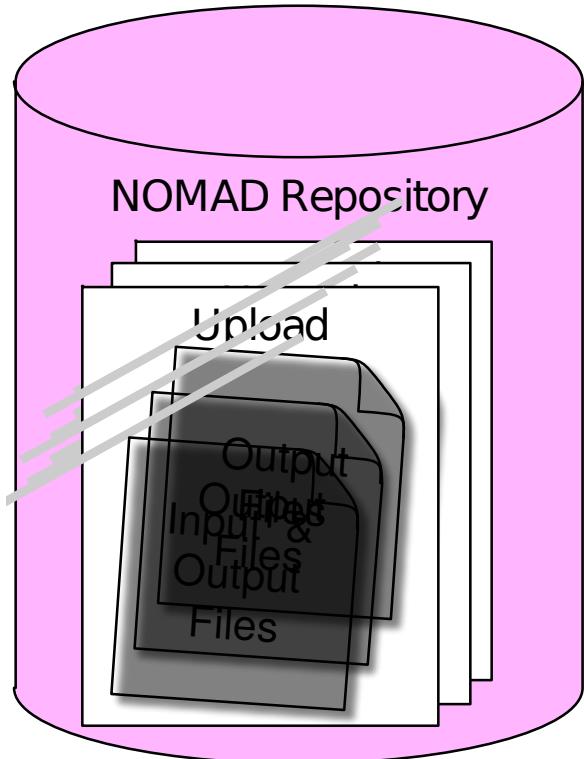
NOVEL MATERIALS DISCOVERY

# NOMAD Archive: How is it built





NOVEL MATERIALS DISCOVERY



# NOMAD Repository

- <http://nomad-repository.eu>



Claudia Draxl  
HUB



Matthias Scheffler  
FHI



Jungho Shin  
HUB



Lorenzo Pardini  
HUB

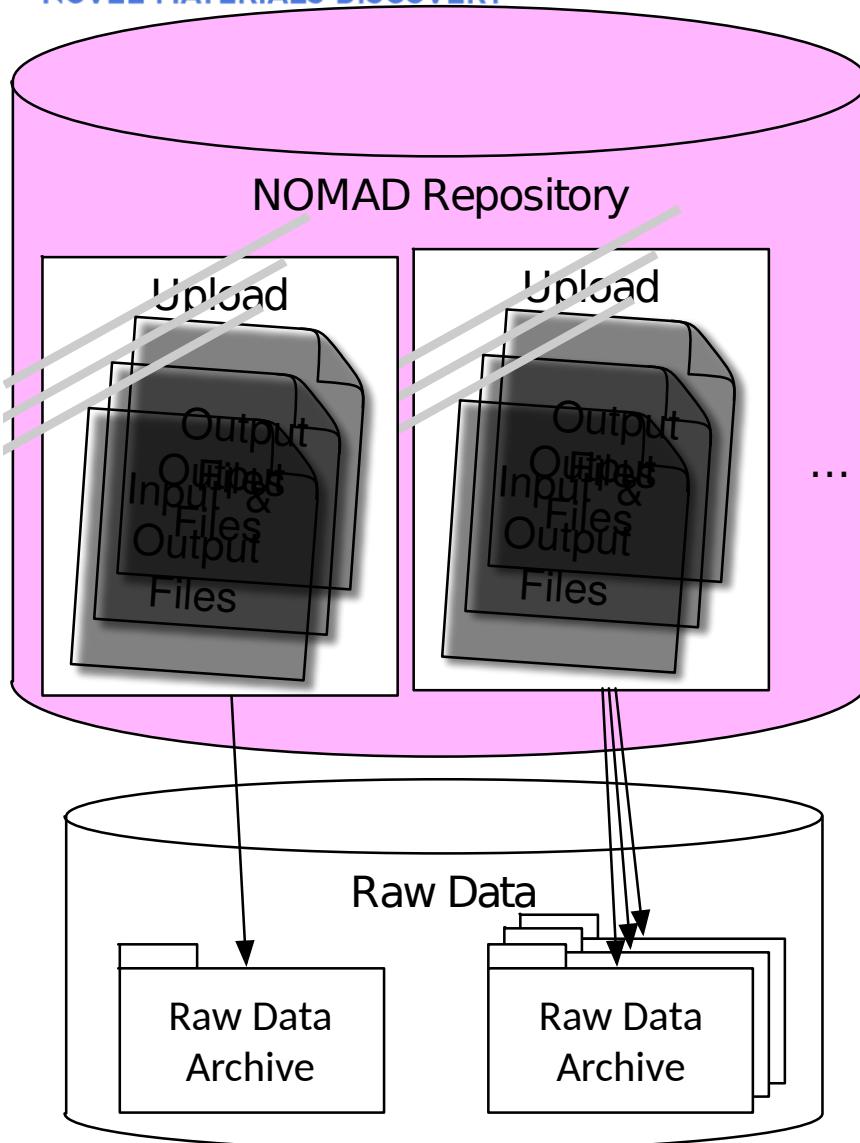


Thomas Zastrow  
MPCDF

- Source of our data
- Established to host organize and share materials data
- Keeps data for at least 10 Years
- Open access and restricted data
  - >3M entries, >2.8M open access
- We use only open access data
- Joint effort by the groups of
  - Matthias Scheffler, FHI Berlin
  - Claudia Draxl, HU Berlin
  - Max Planck Computer & Data Facility (MPCDF), Garching, headed by Stefan Heinzel.



NOVEL MATERIALS DISCOVERY

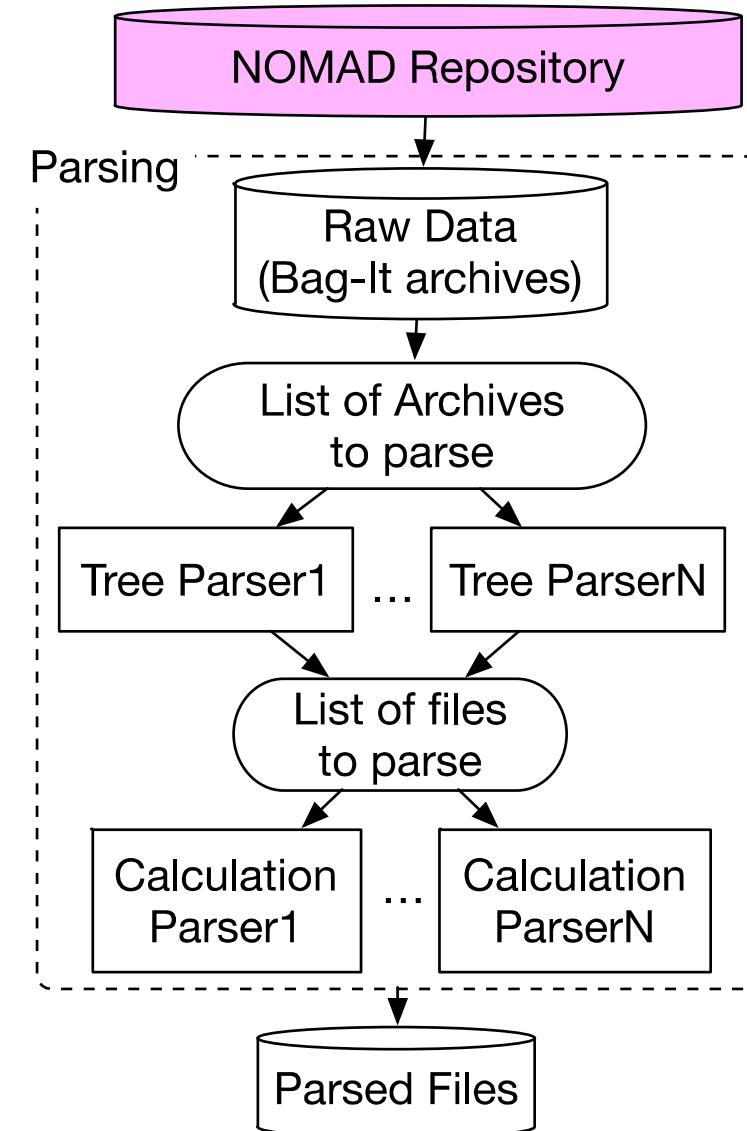


# Raw Data

- **BagIt** format
  - Zip archive
  - Checksums (verifiable integrity)
  - Leading format for digital archiving and exchange
- Limited size (~15GB of uncompressed data) by splitting large uploads
- **Identifier** based on the contents (reproducible)
- Good and important "groundwork" for processing
- Open access: <http://data.nomad-coe.eu/raw-data>



NOVEL MATERIALS DISCOVERY

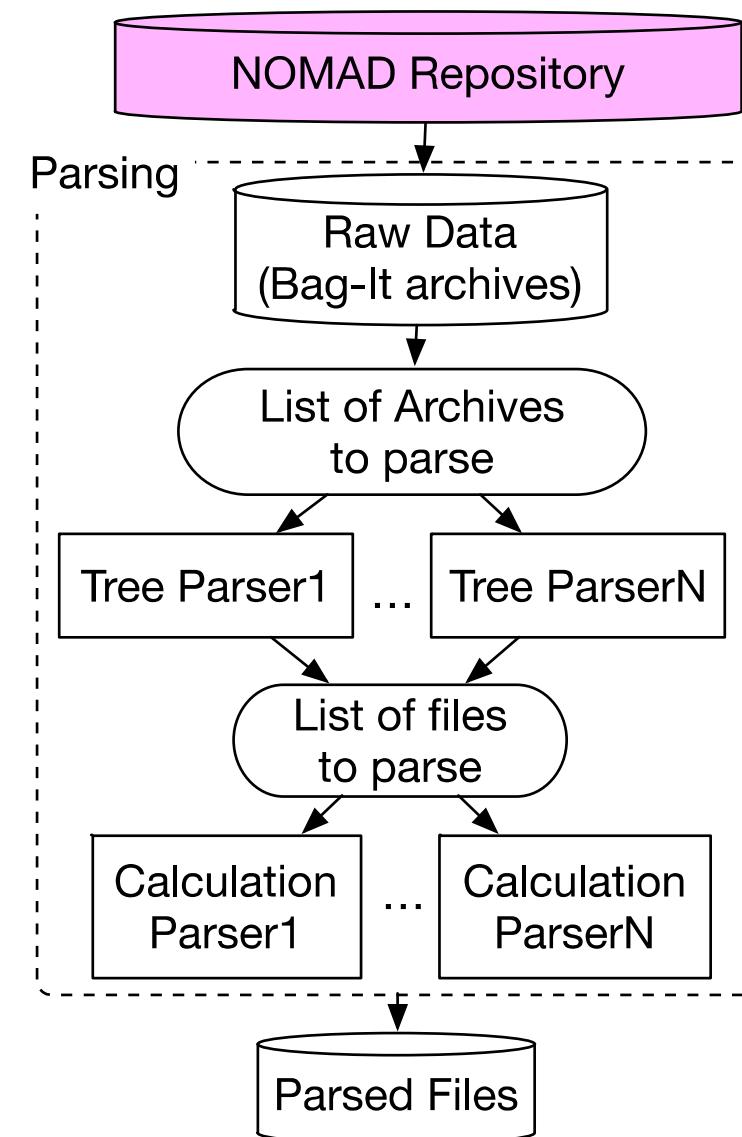


# Parsing

- 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

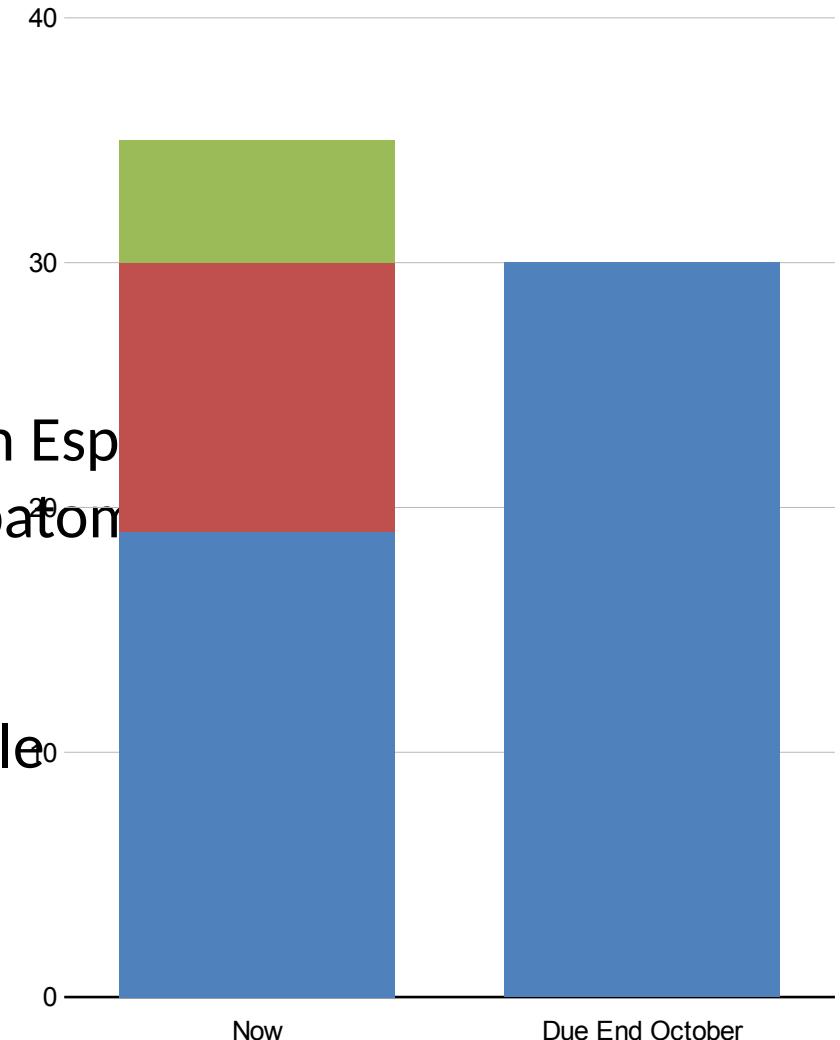


NOVEL MATERIALS DISCOVERY



# Parsers Status

35 Parsers in active development



## Ready:

- abinit
- CASTEP
- cp2k
- CPMD
- FHI-aims
- FLEUR
- Gaussian
- GPAW
- NWChem
- octopus
- onetep
- qBox
- Quantum Espresso
- QUIP /libatoms
- SIESTA
- Smeagol
- turbomole
- VASP
- WIEN2k

# NOMAD

NOVEL MATERIALS  
OPTIMIZATION

Parsing



Mikkel Strange



Micael Oliveira



Martina Stella



Massimo Riello



Wael Chibiani



Franz Knuth



Aliaksei Mazheika

Tree Parsing



Adam Fekete



Adriel Dominiguez



Honghui Shang



Sebastian Alarcón Villaseca



Daria Tomecka



Lorenzo Pardini



Fawzi Mohamed

Calculator  
Parser



Carl Poelking



Rosenndo Valero



Ask Hjorth Larsen



Lauri Himannen



Sami K. Kivistö



Henning Glawe

# Parser Developers

Parsed Files



# Metadata

- **Metadata** is the conceptual model of our data
- Format independent
- Describes both the data and its structure



NOVEL MATERIALS DISCOVERY

# Metadata

*section\_run*

program\_name FHI-aims  
program\_version 081912

**Values:** Data

Structures

and names: Metadata

describes the data  
allows one to express queries  
annotates operations performed on the data



NOVEL MATERIALS DISCOVERY

# Metadata

*section\_run*

program\_name FHI-aims

program\_version 081912

*section\_system*

*section\_method*

*section\_single\_configuration\_calculation*

**Values:** Data

Structures

and names: Metadata

describes the data  
allows one to express queries  
annotates operations performed on the data



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*

## *section\_single\_configuration\_calculation*

**Values:** Data

**Structures**

**and names:** Metadata

describes the data  
allows one to express queries  
annotates operations performed on the data



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*

## *section\_single\_configuration\_calculation*

**Values:** Data

**Structures**

**and names:** Metadata

SI Units:

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

describes the data

allows one to express queries

annotates operations performed on the data



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*

**Values:** Data

**Structures**

**and names:** Metadata

SI Units:

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

describes the data

allows one to express queries

annotates operations performed on the data



NOVEL MATERIALS DISCOVERY

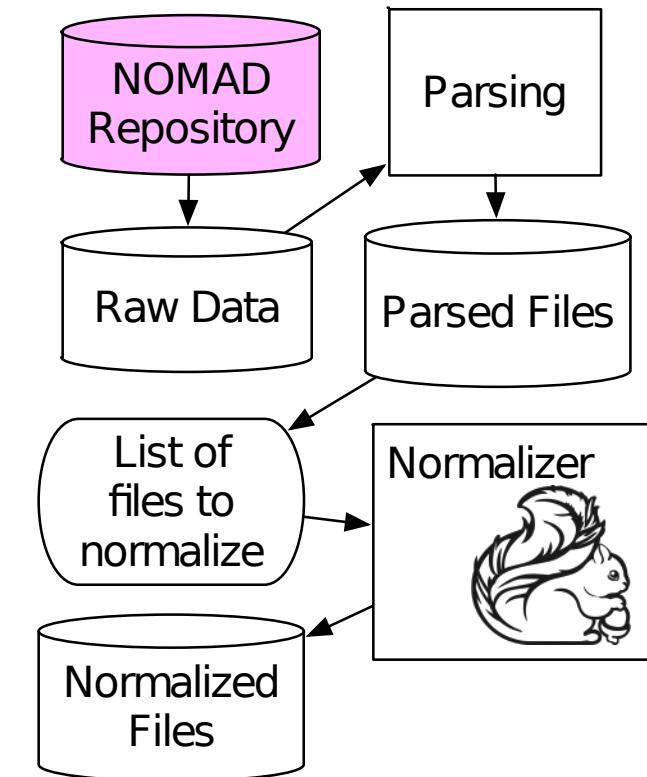
# Normalization

## 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

# Normalized Files



- Parsed and Normalized files are stored using
  - JSON: human readable, nice to use in the web
  - HDF5: efficient binary representation, indexed access
- Raw data identifier used to group calculations
- HDF5 files and raw data available on  
<http://data.nomad-coe.eu/>



NOVEL MATERIALS DISCOVERY

# NOMAD Archive: What is in it





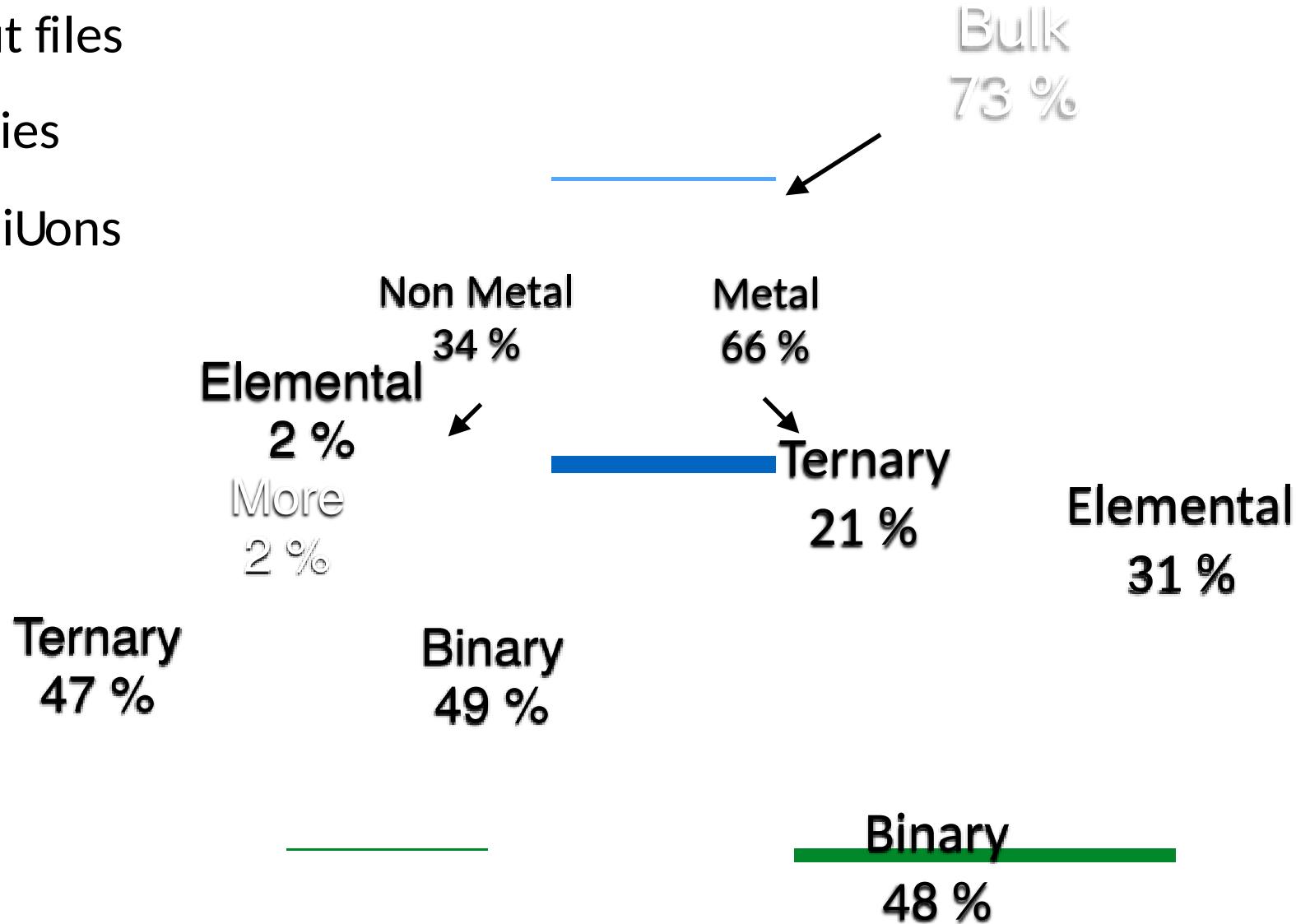
NOVEL MATERIALS DISCOVERY

# **NOMAD Archive: What is in it**

# Molecule

27 %

- >2.6M successfully parsed output files
  - energy of >15M unique geometries
  - >288k different material compositions  
(>41k reduced)
  - >260k band structures
  - >2.4M eigenvalues



# NOMAD Archive: What is in it

Molecule  
27 %

Bulk  
73 %

Atom

33.663

Molecule

5.864.435

Surface

37.632

Bulk

15.939.034

1

100

10.000

1.000.000

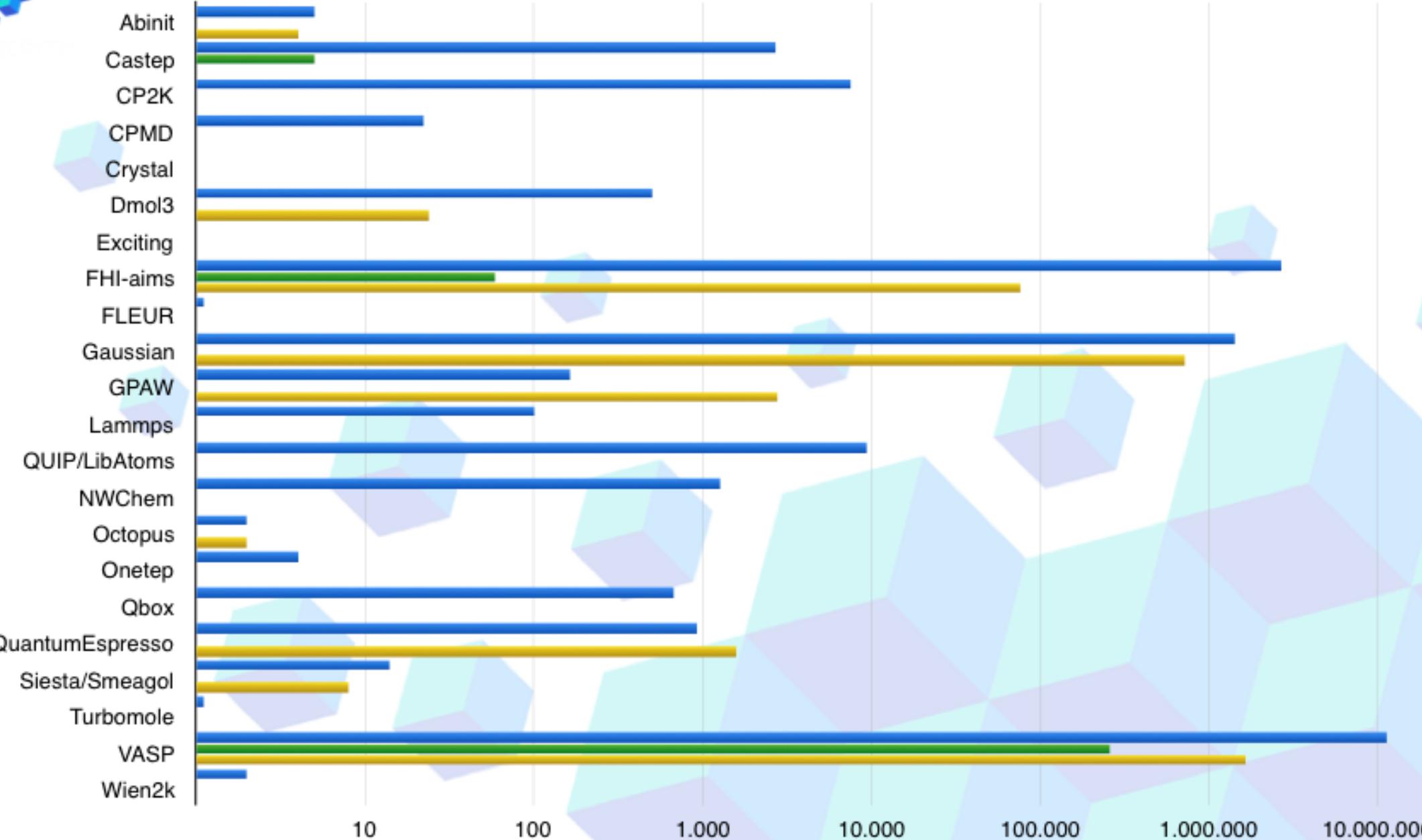
100.000.000

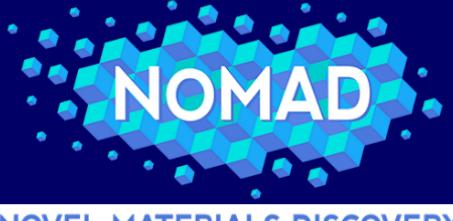


■ Unique Geometries

■ Bands

■ Eigenvalues





NOVEL MATERIALS DISCOVERY

# NOMAD Archive: Conclusions

The core development of the NOMAD Archive is done

- 19 codes are ready and 16 new ones are in development
- An extensive metadata description for *ab initio* calculations
- Standard formats to store the data
- Open access to the archive: <http://data.nomad-coe.eu>

Future

- Parser development to fix errors and support new code versions and codes
- Roughly weekly reparse of data with new parsers
- Extend normalization with new transformations used by the encyclopedia (WP2), visualization (WP3) and data analysis (WP4)
- Improve parser speed
- Extend and complete the metadata (force-field codes, excited states,...)

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 676580.

