



NOVEL MATERIALS DISCOVERY

WP1 Report

Status of the NOMAD Archive

NOMAD SAC Meeting Barcellona 4-5 Oct. 2016

Fawzi Mohamed FHI Berlin

NOMAD

NOVEL MATERIALS DISCOVERY

NOMAD Archive: Goal





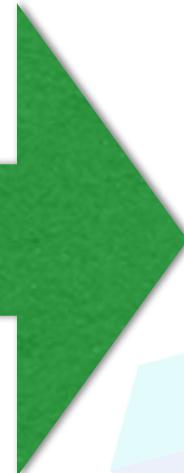
NOVEL MATERIALS DISCOVERY

NOMAD
REPOSITORY

NOMAD Archive: Goal

NOMAD Archive

Make data accessible for
the simulation codes used
in the community



WP2: NOMAD
ENCYCLOPEDIA

WP3: NOMAD
VISUALIZATION

WP4: BIG DATA
ANALYSIS

External users

NOMAD

NOVEL MATERIALS DISCOVERY

NOMAD Archive: What is it



NOVEL MATERIALS DISCOVERY

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



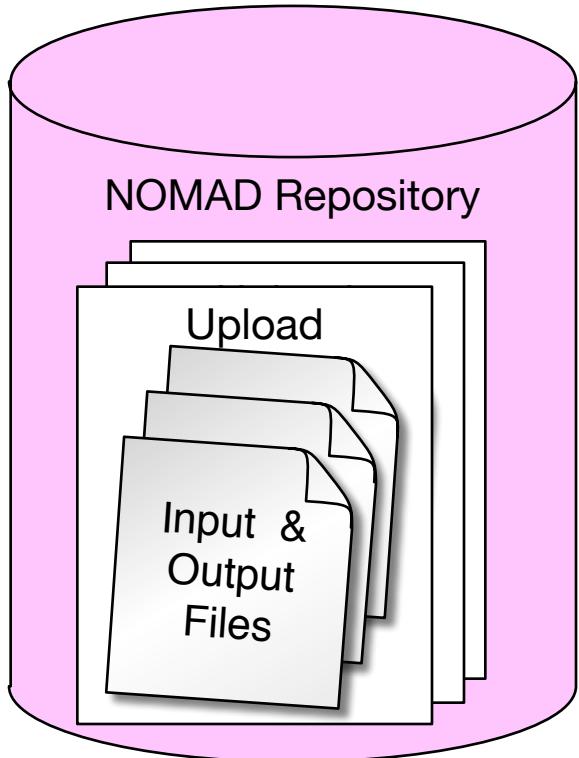
NOVEL MATERIALS DISCOVERY

NOMAD Archive: How is it built





NOVEL MATERIALS DISCOVERY



NOMAD Repository

- <http://nomad-repository.eu>
- Source of our data
- established to host organize and share materials data
- Keeps data for at least 10 Years
- Open access
- 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.



Claudia Draxl
HUB



Matthias Scheffler
FHI



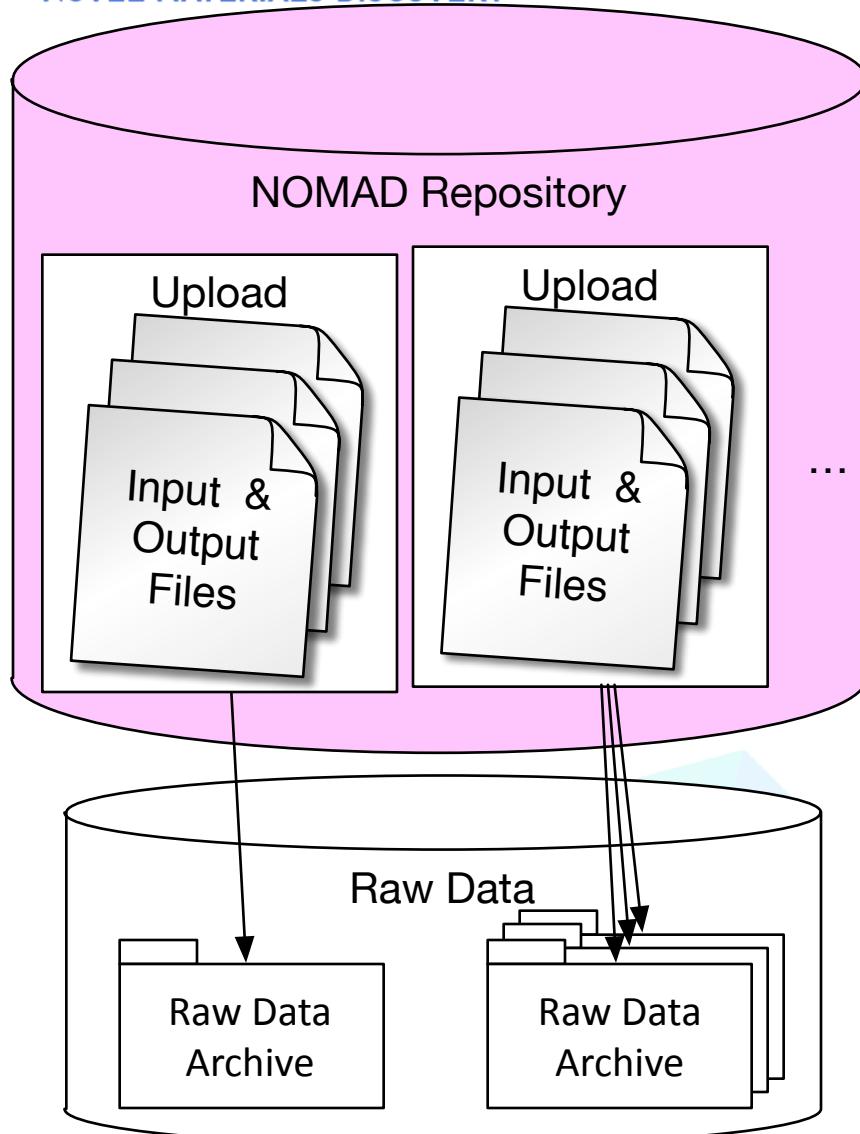
Jungho Shin
HUB



Lorenzo Pardini
HUB

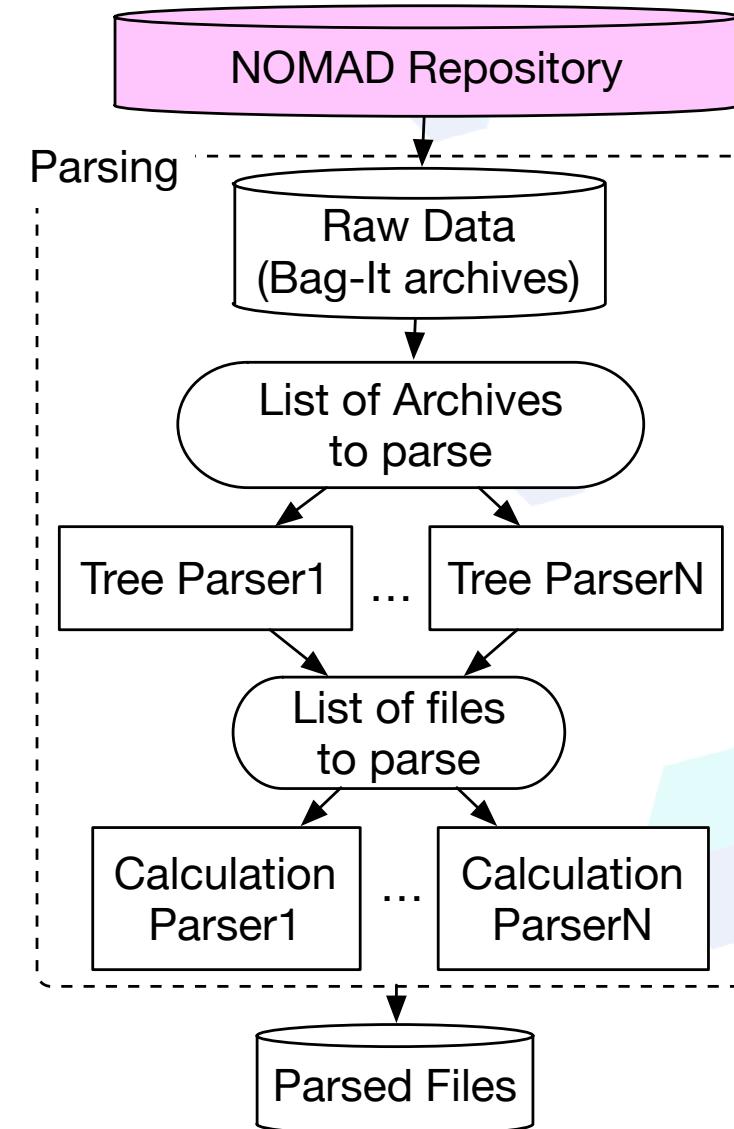


Thomas Zastrow
MPCDF



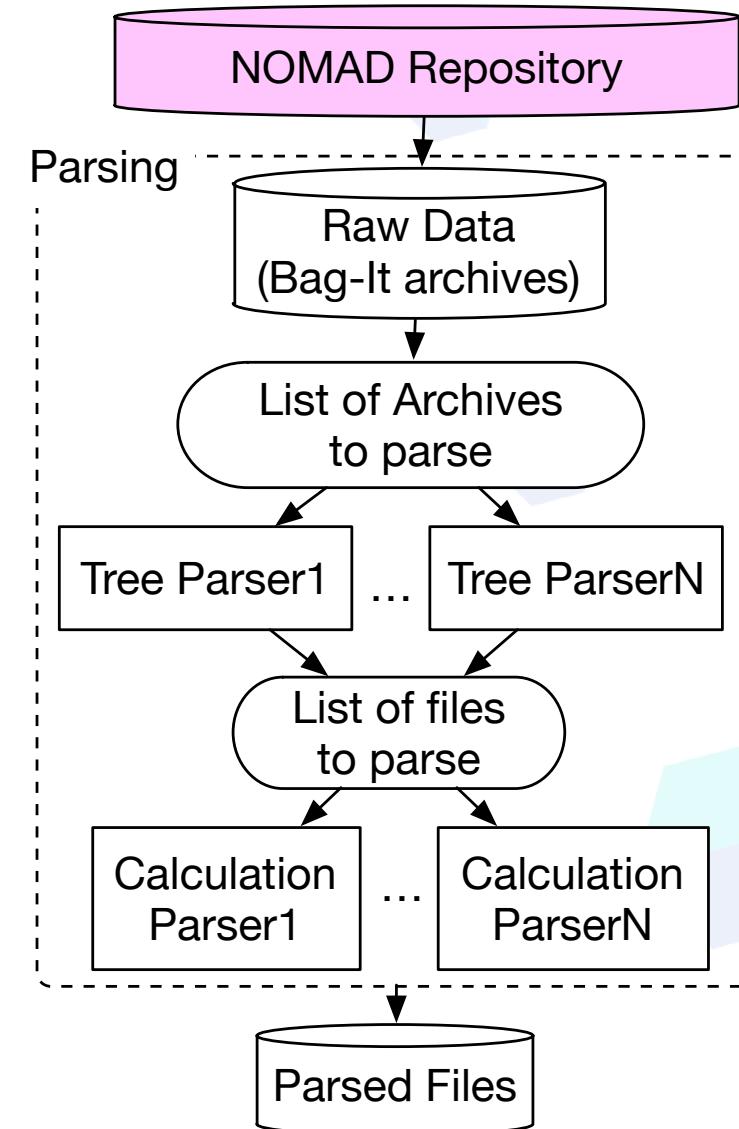
Raw Data

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



Parsing

- **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
- 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*



Parsers Ready

- FHI-aims
- VASP
- Turbomole
- qBox
- ORCA
- exciting
- WIEN2k
- ELK
- Gaussian
- SIESTA
- LAMMPS
- CASTEP
- onetep
- DL_POLY
- QUIP/libatoms/GAP
- CP2K
- Crystal
- CPMD
- NWChem
- Octopus
- Quantum Espresso
- Smeagol
- abinit
- GPAW
- Mopac

NOMAD

NOVEL MATERIALS
OPTIMIZATION

Parsing

Mikkel Strange

Micael Oliveira

Martina Stella

Massimo Riello

Wael Chibiani

Franz Knuth

Aliaksei Mazheika

Andrea Droghetti

Sebastian Alarcón Villaseca

Lorenzo Pardini

Adam Fekete

Honghui Shang

Daria Tomecka

Fawzi Mohamed

Adriel Dominguez

Ask Hjorth Larsen

Sami K. Kivistö

Calculus
Parser

Carl Poelking

Rosenndo Valero

Lauri Himannen

Henning Glawe

Parsed Files

Parser Developers

Metadata

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

Metadata

section_run

program_name

FHI-aims

program_version

081912

Values: Data

Structures

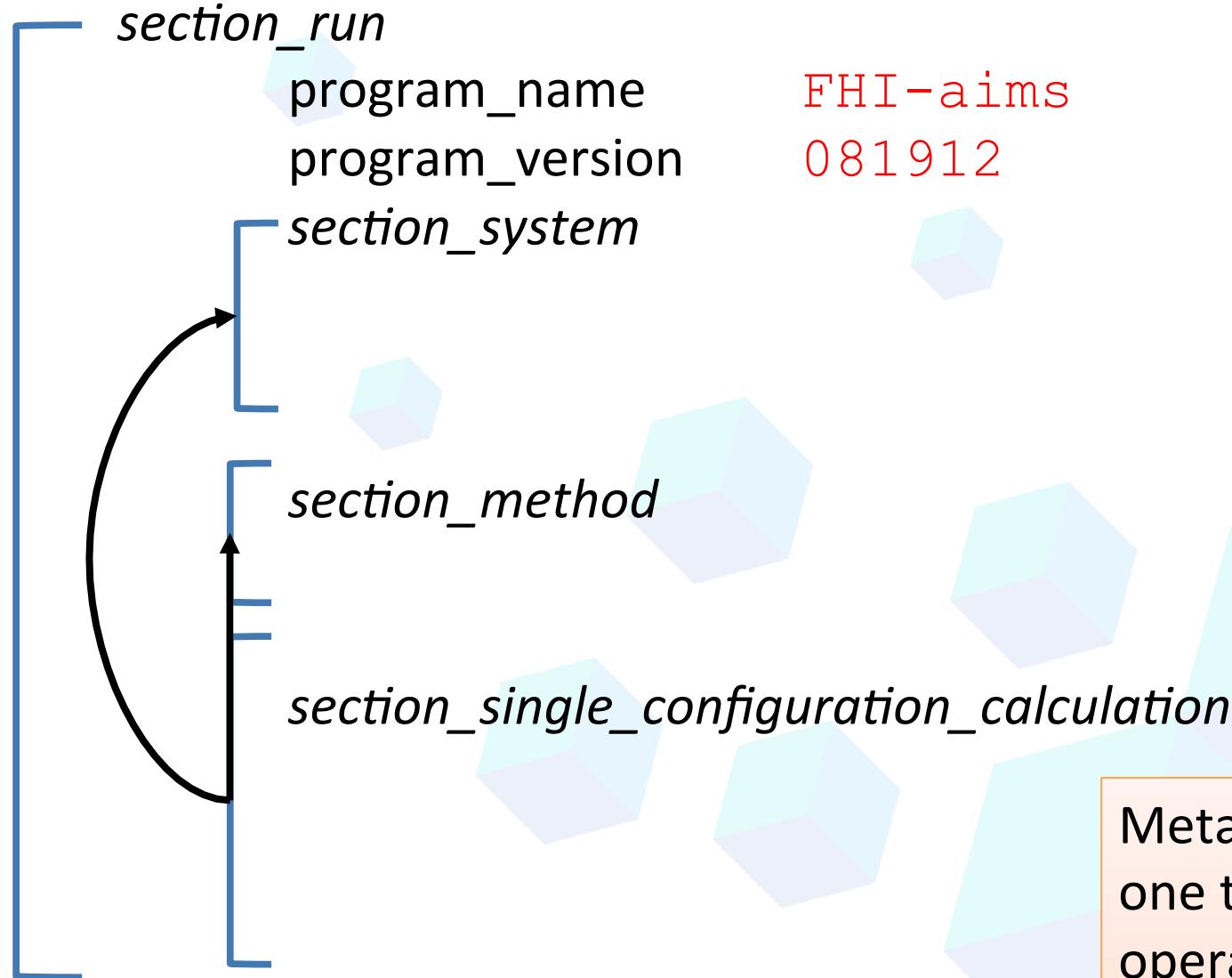
and names: Metadata

Metadata describes the data, and allows one to express queries, and annotate operations performed on the data



NOVEL MATERIALS DISCOVERY

Metadata



Values: Data
Structures
and names: Metadata

Metadata describes the data, and allows one to express queries, and annotate operations performed on the data

Metadata

section_run

program_name	FHI-aims
program_version	081912

section_system

simulation_cell	$[[1.4e-9 \dots]]$
atom_positions	$[[0.0, \dots] \dots]$
atom_labels	$["Cu", \dots]$

section_method

section_single_configuration_calculation

Values: Data
Structures
and names: Metadata

Metadata describes the data, and allows one to express queries, and annotate operations performed on the data

Metadata

section_run

program_name	FHI-aims
program_version	081912

section_system

simulation_cell	$[[1.4e-9 \dots]]$
atom_positions	$[[0.0, \dots] \dots]$
atom_labels	["Cu", ...]

section_method

section_single_configuration_calculation

Values: Data

Structures

and names: Metadata

SI Units:

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

Metadata describes the data, and allows one to express queries, and annotate operations performed on the data

Metadata

<i>section_run</i>	
program_name	FHI-aims
program_version	081912
<i>section_system</i>	
simulation_cell	<code>[[1.4e-9 ...]]</code>
atom_positions	<code>[[0.0, ...], ...]</code>
atom_labels	<code>["Cu", ...]</code>
<i>section_method</i>	
basis_set	fhi_aims_tight
XC_method	DFT_GGA_PBE
<i>section_single_configuration_calculation</i>	

Values: Data

Structures

and names: Metadata

SI Units:

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

Metadata describes the data, and allows one to express queries, and annotate operations performed on the data

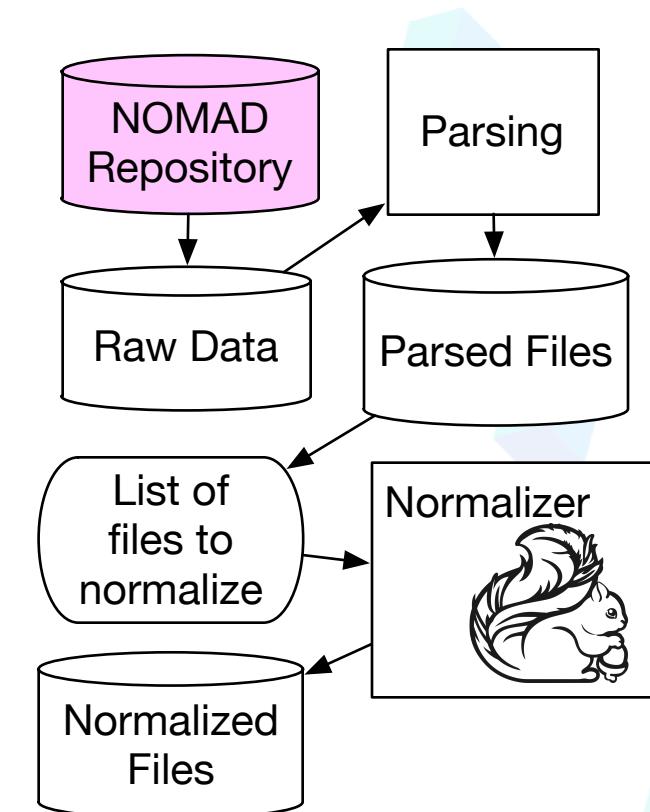
Normalization

Standardization

- Parsers **standardizes** the data, **avoiding to loose information**

Normalization

- The varius users (WP2, WP3, WP4,...) may define **derived quantities** (normalized representations,...) that 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



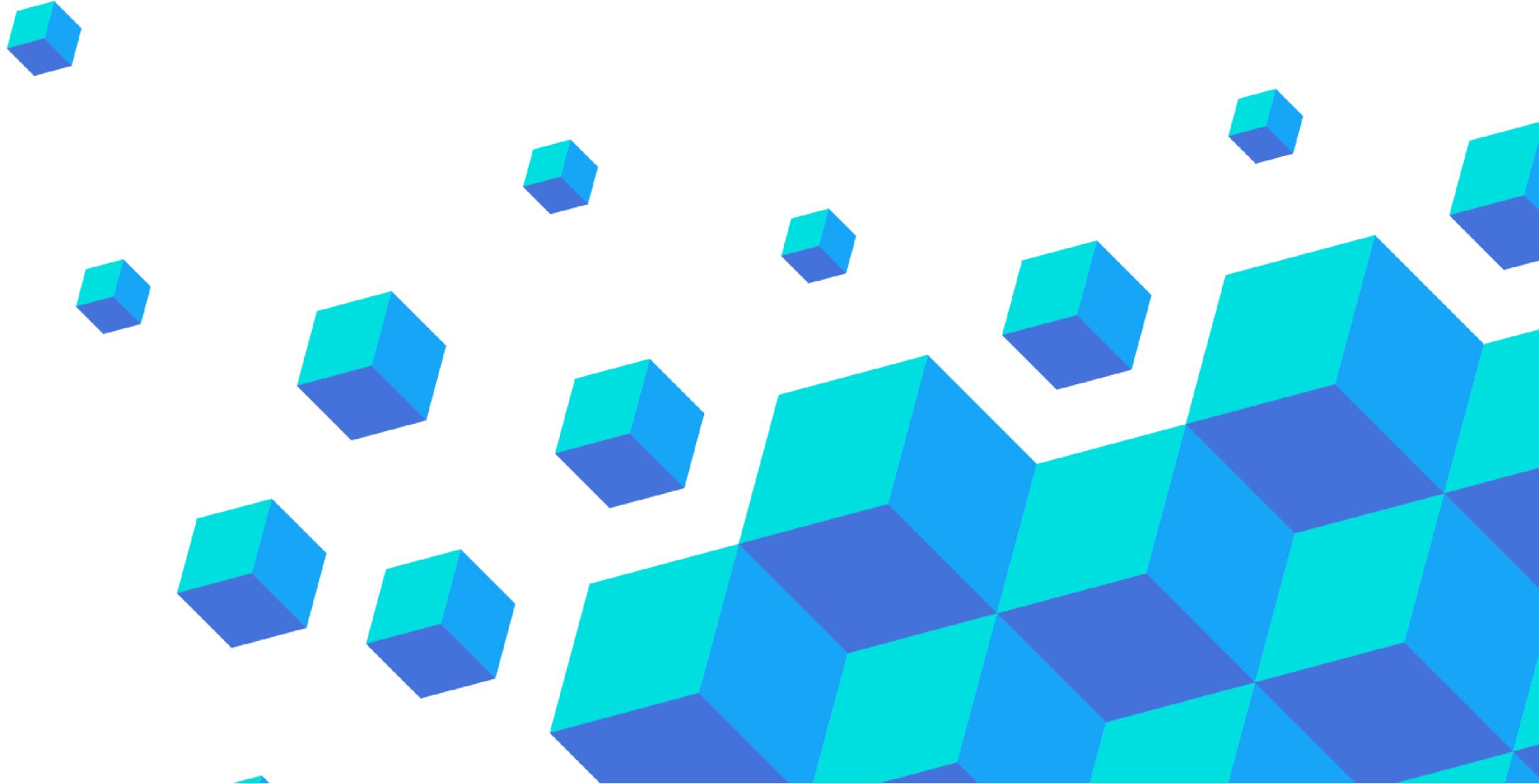
Normalized Files

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



NOVEL MATERIALS DISCOVERY

NOMAD Archive: What is in it



NOMAD Archive: What is in it

Semiconductors or
Insulators

Elemental

2%

Quaternary

1%

Binary
35%

Ternary

62%

Non Metals
59%

Metals
41%

Metals

Elementary

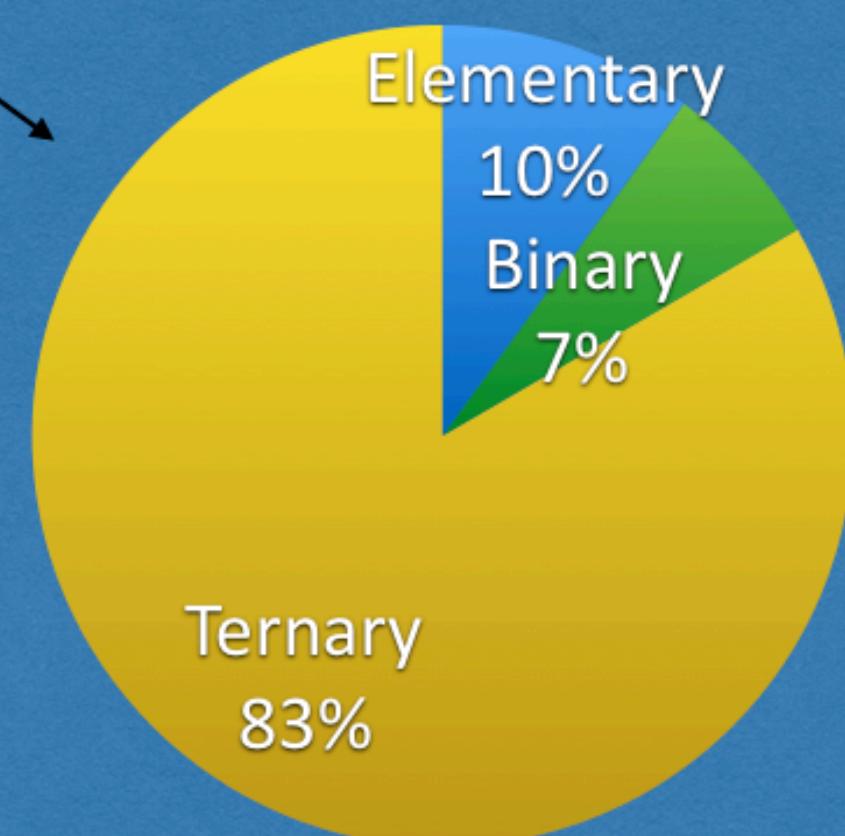
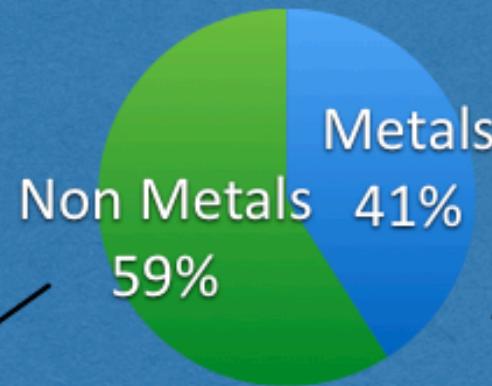
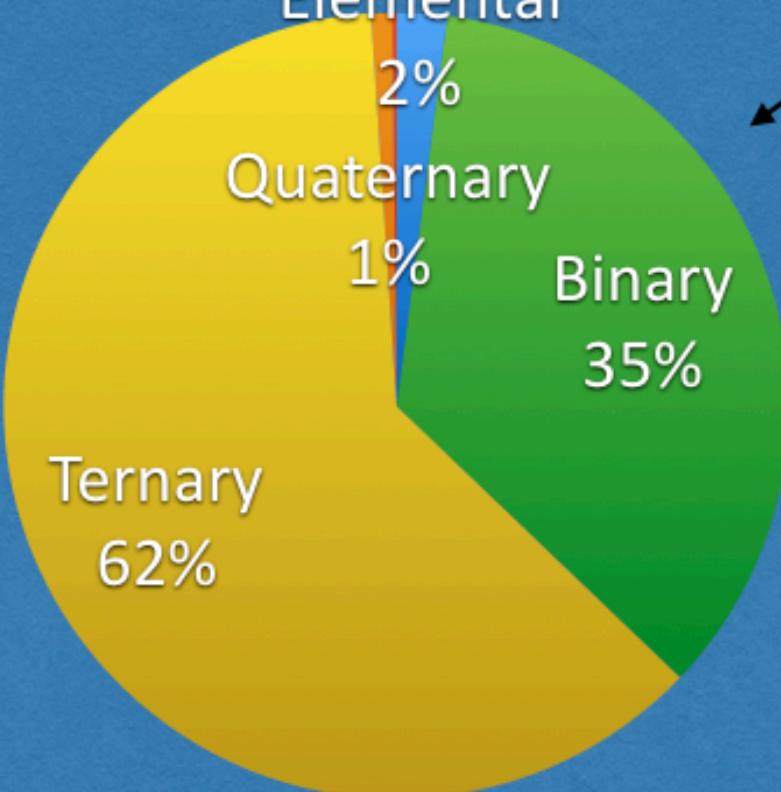
10%

Binary

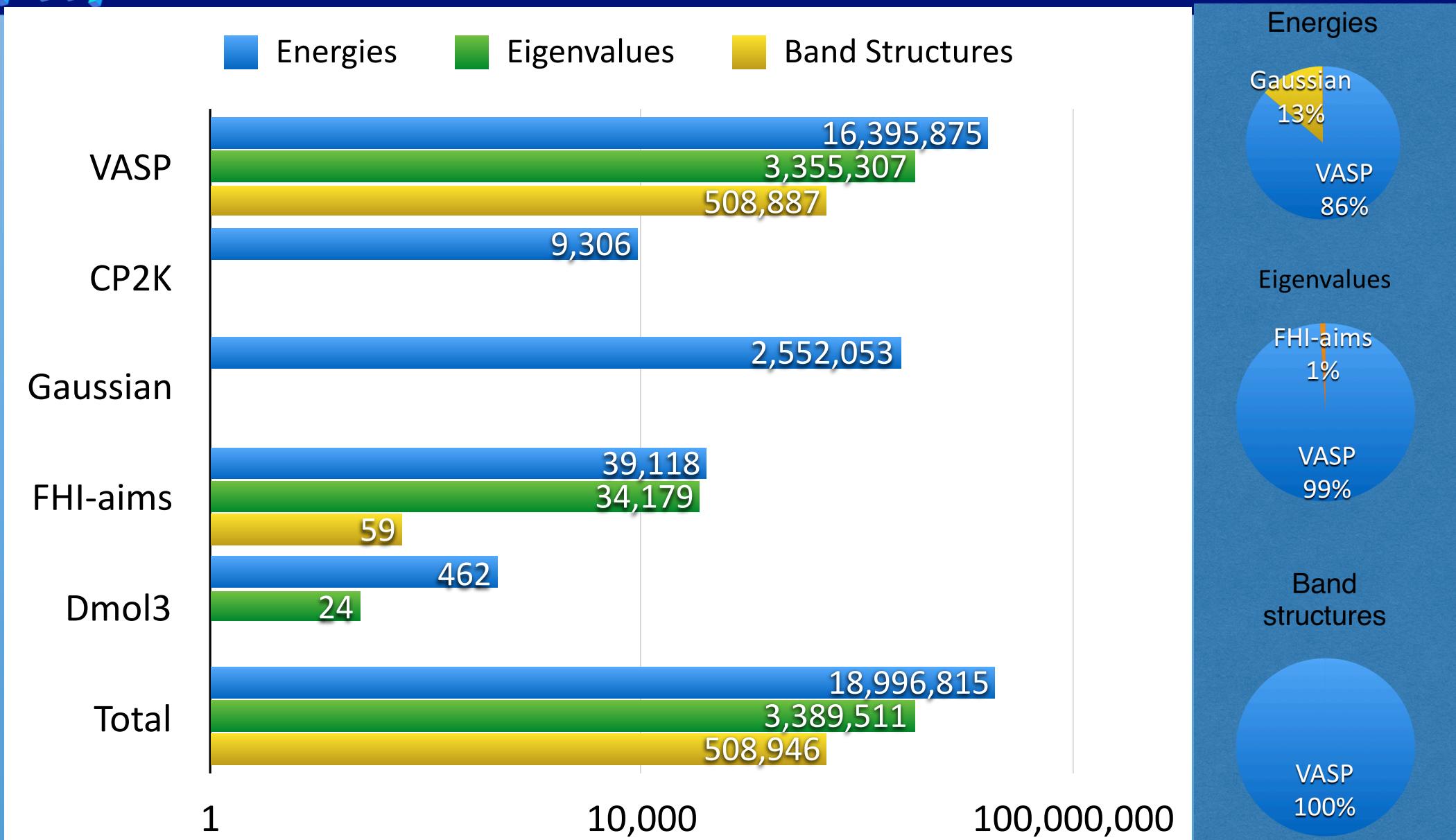
7%

Ternary

83%



NOMAD Archive: What is in it





NOMAD Archive: Conclusions

- The core development of the NOMAD Archive is done,
 - 25 codes are ready and 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 support new code versions and new codes
 - Weekly reparse of data with new parsers
 - Extend Normalization with new transformations used by WP2, WP3 and WP4
 - Improve parser speed
 - Extend metadata to force-field codes