MAY 1 - 5, 2023 - IPAM WORKSHOP III: COMPLEX SCIENTIFIC WORKFLOWS AT EXTREME COMPUTATIONAL SCALES







## Multiscale challenge in materials modeling and experiment



Inherent multiscale character of materials



Multiscale & multiscale simulations often require complex simulation workflows and large-scale simulations

> But: Tools have been independently developed and are not interoperable

## **Challenge: Multiscale and Multiphysics Character**





Real Environments Thermodynamics ➤ T, p, μ Electrochemical

PH, U, E<sub>field</sub>

Mechanical

•••

γσ

### Virtual Materials Design

> Predict response of material to thermodynamic, electrochemical, mechanical environments

Exploring and navigating high-dimensional configuration spaces





Exploring and navigating high-dimensional configuration spaces

> Thermal complexity









Exploring and navigating high-dimensional configuration spaces

Thermal complexity



> Chemical complexity









Exploring and navigating high-dimensional configuration spaces

> Thermal complexity



Chemical complexity



Structural complexity





## Key Challenge – Sampling Huge Configuration Spaces



> Exploration of a huge (8N dimensional) configuration space!

- > In principle well-suited for exascale computing
- > But any brute force approach is unfeasible:
  - > 10 data points on each coordinate
  - ➤ 100 atoms
  - ➤ 10<sup>800</sup> configurations

Dimensionality reduction is mandatory!

Math + ML concepts



# Practical examples of sampling/navigating in high-dimensional configuration spaces

## Ab initio up to the melting temperature



X. Zhang, T. Hickel, B. Grabowski, JN, Comp. Mat. Sci. (2018)

## Ab initio up to the melting temperature





Max-Planck-Institut für Eisenforschung GmbH | Jörg Neugebauer

X. Zhang, T. Hickel, B. Grabowski, JN, Comp. Mat. Sci. (2018)

## Fully autonomous algorithm to determine melting point



5000



Max-Planck-Institut für Eisenforschung GmbH | Jörg Neugebauer

Zhu, Janssen, Ishibashi, Kormann, Grabowski, JN, Comp. Mat. Sci. 187, 12 (2021)

## High-entropy alloy discovery





## Machine learning-enabled high-entro alloy discovery

Ziyuan Rao<sup>1</sup>, Po-Yen Tung<sup>1,2</sup>, Ruiwen Xie<sup>3</sup>, Ye Wei<sup>1</sup>\*, Hongbin Zhang<sup>3</sup>, Alberto Ferrari<sup>4</sup>, T.P.C. Klaver<sup>4</sup>, Fritz Körmann<sup>1,4</sup>, Prithiv Thoudden Sukumar<sup>1</sup>, Alisson Kwiatkowski da Silva<sup>1</sup>, Yao Chen<sup>1,5</sup>, Zhiming Li<sup>1,6</sup>, Dirk Ponge<sup>1</sup>, Jörg Neugebauer<sup>1</sup>, Oliver Gutfleisch<sup>1,3</sup>, Stefan Bauer<sup>7</sup>, Dierk Raabe<sup>1\*</sup>



## Many of the recent breakthroughs in materials science simulations rely on complex simulation protocols (workflows)



## **Workflows: Fundamentals**





## Workflows: Challenges





is mandatory

## Libraries that successfully address some of these issues



>>> import numpy as np
>>> from scipy.optimize import minimize

```
>>> def rosen(x):
... """The Rosenbrock function"""
... return sum(100.0*(x[1:]-x[:-1]**2.0)**2.0 + (1-x[:-1])**2.0)
```

```
>>> print(res.x)
[1. 1. 1. 1. 1.]
```

## Libraries that successfully address some of these issues



>>> import numpy as np
>>> from scipy.optimize import minimize



```
>>> print(res.x)
[1. 1. 1. 1. 1.]
```

Fully embedded in python!

## **Translation to materials science simulations**







>>> print(res.x) [1. 1. 1. 1. 1.]

Fully embedded in python!

## Translation to materials science simulations

Features beyond NumPy/SciPy

- Automatic upload of input + selected output to database/storage
  - Serialization
- Upscaling to HPC (exascale)
  - Job managment
- Full integration into Jupyter
  - > e.g. visualization of complex data structures such as atomic trajectories



### Fully embedded in python!



## Integrated development environment (IDE) for workflows





Max-Planck-Institut für Eisenforschung GmbH | Jörg Neugebauer

www.pyiron.org

## Key concept: pyiron objects

## Identical generic look and feel



job.structure job.input job.output.energy\_tot job.run() job.structure job.input job.output.energy\_tot job.run() job.structure job.input job.output.energy\_tot job.run()































Structure

Relaxation

Solid

Melting

Liquid

Combine

Interface



### Create structure

fcc Al with vacancy

```
Cu = pr.create.structure.bulk('Cu', cubic=True).repeat(3)
Cu[19] = 'Al'
Cu.plot3d(particle_size=2)
```









### Perform MD simulations

- job = pr.create.job.Lammps(job\_name='MD\_700K')
  job.structure = Cu
  job.list\_potentials()[:10]
- ['1999--Liu-X-Y--Al-Cu--LAMMPS--ipr1', '2011--Apostol-F--Al-Cu--LAMMPS--ipr1', '2012--Jelinek-B--Al-Si-Mg-Cu-Fe--LAMMPS--ipr2', '2016--Zhou-X-W--Al-Cu--LAMMPS--ipr2', '2018--Zhou-X-W--Al-Cu-H--LAMMPS--ipr1', '2022--Mahata-A--Al-Cu--LAMMPS--ipr1', 'EAM\_Dynamo\_CaiYe\_1996\_AlCu\_\_M0\_942551040047\_005', 'EAM\_Dynamo\_LiuLiuBorucki\_1999\_AlCu\_\_M0\_020851069572\_000', 'EMT\_Asap\_Standard\_JacobsenStoltzeNorskov\_1996\_AlAgAuCuNiPdPt\_\_M0\_115316750986\_001', 'MEAM\_LAMMPS\_JelinekGrohHorstemeyer\_2012\_AlSiMgCuFe\_\_M0\_262519520678\_000']







### Perform MD simulations

- job = pr.create.job.Lammps(job\_name='MD\_700K')
  job.structure = Cu
  job.list\_potentials()[:10]
- ['1999--Liu-X-Y--Al-Cu--LAMMPS--ipr1', '2011--Apostol-F--Al-Cu--LAMMPS--ipr1', '2012--Jelinek-B--Al-Si-Mg-Cu-Fe--LAMMPS--ipr2', '2016--Zhou-X-W--Al-Cu--LAMMPS--ipr1', '2018--Zhou-X-W--Al-Cu-H--LAMMPS--ipr1', '2022--Mahata-A--Al-Cu--LAMMPS--ipr1', 'EAM\_Dynamo\_CaiYe\_1996\_AlCu\_\_M0\_942551040047\_005', 'EAM\_Dynamo\_LiuLiuBorucki\_1999\_AlCu\_\_M0\_020851069572\_000', 'EAM\_Dynamo\_LiuLiuBorucki\_1999\_AlCu\_\_M0\_020851069572\_000', 'EAM\_Dynamo\_LiuLiuBorucki\_1999\_AlCu\_\_M0\_020851069572\_000', 'EAM\_Dynamo\_LiuLiuBorucki\_1999\_AlCu\_\_M0\_020851069572\_000', 'EAM\_LAMMPS\_JelinekGrohHorstemeyer\_2012\_AlSiMgCuFe\_\_M0\_262519520678\_000'] job.potential = '2022--Mahata-A--Al-Cu--LAMMPS--ipr1' job.calc\_md(temperature=700, n\_ionic\_steps=10000)

```
job.run(run_again=True)
```





Structure
Relaxation
Solid
Melting
Liquid
Combine
Interface

Direct access of all simulation data from database Load job from database and analyse it

	Parameter	Value
0	units	metal
1	dimension	3
2	boundary	qqq
3	atom_style	atomic
4	read_data	structure.inp
5	include	potential.inp
6	fixensemble	all nvt temp 700.0 700.0 0.1
7	variabledumptime	equal 100
8	variablethermotime	equal 100
9	timestep	0.001
10	velocity	all create 1400.0 53471 dist gaussian
11	dump1	all custom \${dumptime} dump.out id type xsu ysu zsu fx fy fz vx vy vz
12	dump_modify1	sort id format line "%d %d %20.15g
3	thermo_style	custom step temp pe etotal pxx pxy pxz pyy pyz pzz vol

format float %20.15g

14

thermo\_modify













### Jobs that are already in the database are automatically loaded rather than rerun

### Run over series of temperatures

```
for T in np.arange(700, 1800, 100):
    job = pr.create.job.Lammps(job_name=f'MD_{T}K')
    job.structure = Cu
```

```
job.potential = '2022--Mahata-A--Al-Cu--LAMMPS--ipr1'
job.calc_md(temperature=T, n_ionic_steps=10000, pressure=0)
```

job**.run**()

2023	3-04-	-30	19:28	:51,	123 -	pyir	on_log -	WARN	ING	– Th	e job	MD_700	< is	being	loaded	instead	of	running.
The	job	MD_	800K	was s	saved	and	received	the	ID:	282								
The	job	MD_	900K	was s	saved	and	received	the	ID:	283								
The	job	MD_	1000K	was	saved	and	receive	d the	ID:	284								
The	job	MD_	_1100K	was	saved	and	receive	d the	ID:	285								
The	job	MD_	_1200K	was	saved	and	receive	d the	ID:	286								
The	job	MD_	_1300K	was	saved	and	receive	d the	ID:	287								



Structure

Relaxation



### Analyse data using pyiron tables

```
def get_temperature(job):
    return int(job.job_name.split('_')[1][:-1])
```

```
def get_displacement(job):
    pos = job['output/generic/unwrapped_positions']
    return np.mean(np.linalg.norm(pos[-1] - pos[0], axis=-1))
```

```
table = pr.create.table(delete_existing_job=True)
table.add['temperature'] = get_temperature
table.add['displacement'] = get_displacement
table.add.get_volume
table.add.get_energy_tot_per_atom
table.add.get_job_name
table.run()
```

The job table was saved and received the ID: 290

Processing jobs: 100%



12/12 [00:00<00:00, 261.72it/s]



pyiron tables provide a powerful tool to map/reduce complex data into pandas Dataframes







pyiron tables provide a powerful tool to map/reduce complex data into pandas Dataframes df = table.get\_dataframe().sort\_values(by='temperature')
df

	job_id	job_name	energy_tot	volume	temperature	displacement
0	278	MD_700K	-3.353856	1270.238787	700	0.181369
4	282	MD_800K	-3.328802	1328.153374	800	0.247824
5	283	MD_900K	-3.256737	1352.336401	900	0.321491
6	284	MD_1000K	-3.286276	1334.194147	1000	0.264272
7	285	MD_1100K	-3.233808	1353.052931	1100	0.325095
8	286	MD_1200K	-3.241711	1346.373954	1200	0.310809
9	287	MD_1300K	-3.218381	1358.605471	1300	0.338850
1	279	MD_1400K	-3.203465	1374.474437	1400	0.362491
10	288	MD_1500K	-3.110352	1374.445628	1500	0.406013
3	281	MD_1600K	-3.075405	1382.937737	1600	0.646145















### Analyse structure

job\_MD.get\_structure(-1).plot3d(particle\_size=2)







#### Pair-correlation function



Max-Planck-Institut für Eisenforschung GmbH | Jörg Neugebauer

etc.

## Translation to materials science simulations

Features beyond NumPy/SciPy

- Automatic upload of input + selected output to database/storage
  - Serialization
- Upscaling to HPC (exascale)
  - Job managment
- Full integration into Jupyter
  - > e.g. visualization of complex data structures such as atomic trajectories



Pyiron provides all these features!



HP(

(inp, wf, out)





Structure Relaxation Solid Melting Liquid Combine Interface

```
job.potential = 'CuAl_lammps_eam'
job.calc_md(temperature=600, pressure=0, n_ionic_steps=1000)
job.server.queue_view
# job.run()
```

	maximum cores	minimum cores	run time limit
impi_hy*	1280	40	259200
impi_hydra	4240	20	259200
impi_hydra_cmfe.*	1280	40	259200
impi_hydra_small	40	1	604800

```
job.potential = 'CuAl_lammps_eam'
job.calc_md(temperature=600, pressure=0, n_ionic_steps=1000000)
job.server.core = 40
job.server.queue = 'impi_hydra_small'
job.run()
```

## Translation to materials science simulations

Features beyond NumPy/SciPy

- Automatic upload of input + selected output to database/storage
  - Serialization
- Upscaling to HPC (exascale)
  - Job managment
- Full integration into Jupyter
  - > e.g. visualization of complex data structures such as atomic trajectories



Pyiron provides all these features!





### Run over several atomistic simulations Structure pr = Project('benchmark') Al = pr.create.structure.bulk('Al', cubic=True).repeat(3) Relaxation for job\_type in ['Lammps', 'Vasp', 'Sphinx']: job = pr.create\_job(job\_type=job\_type, job\_name=job\_type) job.structure = Al job.run() Solid Al = job.get\_structure(-1) # get last structure Melting Liquid Combine Interface

## Fully interoperable via generic input/output

## Fully autonomous algorithm to determine melting point



5000



Max-Planck-Institut für Eisenforschung GmbH | Jörg Neugebauer

Zhu, Janssen, Ishibashi, Kormann, Grabowski, JN, Comp. Mat. Sci. 187, 12 (2021)

## Sociological challenges

### How to successfully drive transfer?

- user-friendliness
- bottom-up development
- developers must be active users
- internal & external workshops
- hackathons
- documentation & examples
- ➤ 'killer' applications
- scientific breakthroughs



Transition coordinate

## pyiron as platform for internal & external collaborations



### Hands-on interactive workshops (examples)

From Atomistics to Phase Diagrams

Constructing ML Potentials **Talk by R. Drautz** ADIS 2021

> 200 participants

### Teaching platform (example)

Lecture: Modelling and Engineering of Nanoscale Materials (Univ. Gent)

Quote:

"pyiron changed our exercise lessons drastically, by allowing more time to tackle advanced and complex problems instead of dealing with code technicalities."

– Sander Borgmans (Tutor)

### Hackathons & Journal Club

- Weekly 2-3 hour events
- > 5-10 people
- Interactive prototyping & development of new concepts and tools

Interactive seminars

### **Publications**

Read papers and run published notebooks using pyiron and mybinder

A fully automated approach to calculate the melting temperature of elemental crystals





## **Nucleus for large-scale networks: MaterialDigital**





Initiative to digitize materials: particular strategic importance for Germany as a business location

### The project goals: What an industrially relevant material data space must fulfill

**Sovereignty:** Security, ownership and access rights to the data, taking into account the interests of the data creators, are at the forefront of all activities.

**Reproducibility:** The standardized description of data generation. This must be transparent, traceable and repeatable.



**Accessibility**: The uniformly and comprehensively described data must be retrievable in order to avoid redundancies in research.



Adaptability: Flexible data logic that is always aligned with new research and development findings.

Curation: Clearly defined specifications for the data enable its quality assurance.













## **Nucleus for large-scale networks: MaterialDigital**





Initiative to digitize materials: particular strategic importance for Germany as a business location

### The vision: A decentralized data room as a unity of data and data processing



**Generation of material data**: When data is generated, it is stored according to unified schemas.



**Decentralized storage**: The large amounts of data remain at the place of their creation. No data provider cedes its control.



**Integrated analysis:** Linked software environments allow standardized data access and processing.



**Remote access:** The network architecture allows external partners to authorize access to the local environment.











Max-Planck-Institut für Eisenforschung GmbH | Jörg Neugebauer

### Automated workflows





## Nucleus for large-scale networks: NFDI-MatWerk





## **NFDI-MatWerk**

On the route towards a *National Research Data Infrastructure* for Materials Science & Engineering.

> Workflows as central part in materials science & engineering!



www.nfdi-matwerk.de

## Automating simulation life cycle by pyiron



## **Next steps**



### IronFlow - Jupyter-based visual scripting gui for running pyiron workflow graphs

- Note: under active development
   not yet production ready
- ➤ Features
  - Ontologic typing
  - Fully integrated in Jupyter
  - Batch jobs (e.g.for ML)
  - New nodes can be easily written and added
- Workflow programming and application for everybody



#### Max-Planck-Institut für Eisenforschung GmbH | Jörg Neugebauer

### https://github.com/pyiron/ironflow

## Node-based pyiron architecture





## **Application: Materials design**







### METALLURGY

Rao et al., Science 378, 78-85 (2022)



## Machine learning-enabled high-entropy alloy discovery

Ziyuan Rao<sup>1</sup>, Po-Yen Tung<sup>1,2</sup>, Ruiwen Xie<sup>3</sup>, Ye Wei<sup>1</sup>\*, Hongbin Zhang<sup>3</sup>, Alberto Ferrari<sup>4</sup>, T.P.C. Klaver<sup>4</sup>, Fritz Körmann<sup>1,4</sup>, Prithiv Thoudden Sukumar<sup>1</sup>, Alisson Kwiatkowski da Silva<sup>1</sup>, Yao Chen<sup>1,5</sup>, Zhiming Li<sup>1,6</sup>, Dirk Ponge<sup>1</sup>, Jörg Neugebauer<sup>1</sup>, Oliver Gutfleisch<sup>1,3</sup>, Stefan Bauer<sup>7</sup>, Dierk Raabe<sup>1\*</sup>

Identify from millions of possible compositionally complex alloy compositions those with the lowest/vanishing thermal expansion coefficient (Invar alloys)

Original Invar Alloy (FeNi)

**Example: Alloy Discovery** 

Aim:



### Challenges:

- Sharp local minima
- Interplay between magnetism & phonons
- > 5 elements

## **ML Architecture**





Rao et al., Science 378, 78–85 (2022)

## **Importance of Physics-Informed Descriptors**



### Physics descriptors: $\omega_s$ - magnetostriction $T_c$ - Curie temperature):

### Inclusion of DFT computed descriptors significantly improves ML model!

Rao et al., Science 378, 78–85 (2022)



## **Dimensionality Reduction**





GMM-MCMC sampling

WAE latent space

Islands in latent space reflect compositional differences and provide treasure map for HEA-based Invar alloys!

GMM - Gaussian mixture modelMCMC - Markov chain Monte Carlo samplingWAE - Wasserstein autoencoder architecture

Rao et al., Science 378, 78-85 (2022)

## **Discovery of Invar Alloys**





Experiment

Comparison

Design of high-entropy Invar alloys using very sparse experimental data by employing DFT descriptors based on physical models!

Max-Planck-Institut für Eisenforschung GmbH | Jörg Neugebauer

Rao et al., Science 378, 78–85 (2022)

## Conclusions

Advanced thermodynamic and ML approaches together with physics-based descriptors provide powerful tools to explore and utilize high-dimensional configuration spaces

Workflows/simulation protocols to run these approaches become exceedingly complex

 $\rightarrow$  Our approach: **pyiron** as materials IDE

Enable computational materials design in high-dimensional configuration spaces

Conventior

allovs

0.2 0.4 0.6

0.0

TEC (10°/K)

15











## Thank you for your attention!





### Jan Janssen



Developer Team













DFG

Deutsche



SFB

1394



MATERIALD1G1TAL

RUHR EXPLORES SOLVATION

Structural and

Complexity

**Chemical Atomic** 





Federal Ministry of Education and Research

**CLUSTER OF EXCELLENCE - EXC 1069** 

Alexander von Humboldt Stiftung/Foundation