

# Challenges and Opportunities in using Workflow Technology for Reproducible and Reusable Simulation Protocols

INSTITUTE OF NANOTECHNOLOGY



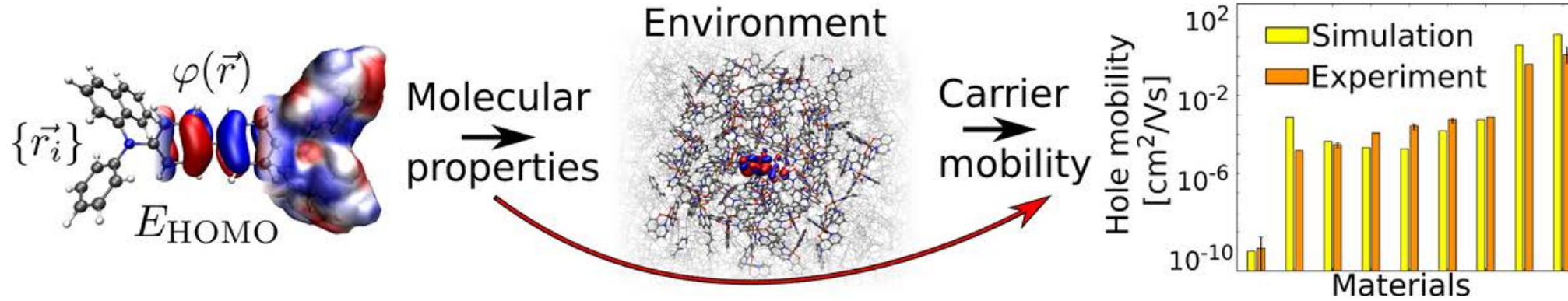
# Challenges and Opportunities of Workflow Technology

- **Multiscale Materials Modelling and Virtual Design @KIT**
- **Challenges in Materials Modelling**
- **Opportunities of Workflow Technology**
- **Challenges of Workflow Technology**

# Challenges and Opportunities of Workflow Technology

- **Multiscale Materials Modelling and Virtual Design @KIT**
- Challenges in Materials Modelling
- Opportunities of Workflow Technology
- Challenges of Workflow Technology

# Multiscale Materials Modelling and Virtual Design AG-Wenzel (INT)



- Development and application of methods for multi-scale simulations of nanoscale materials and devices

→ materials design and discovery

## nanoscale electronics

- single molecule electronics
- organic electronics

## carbon based systems

- graphene and carbon nanotubes

## Workflows

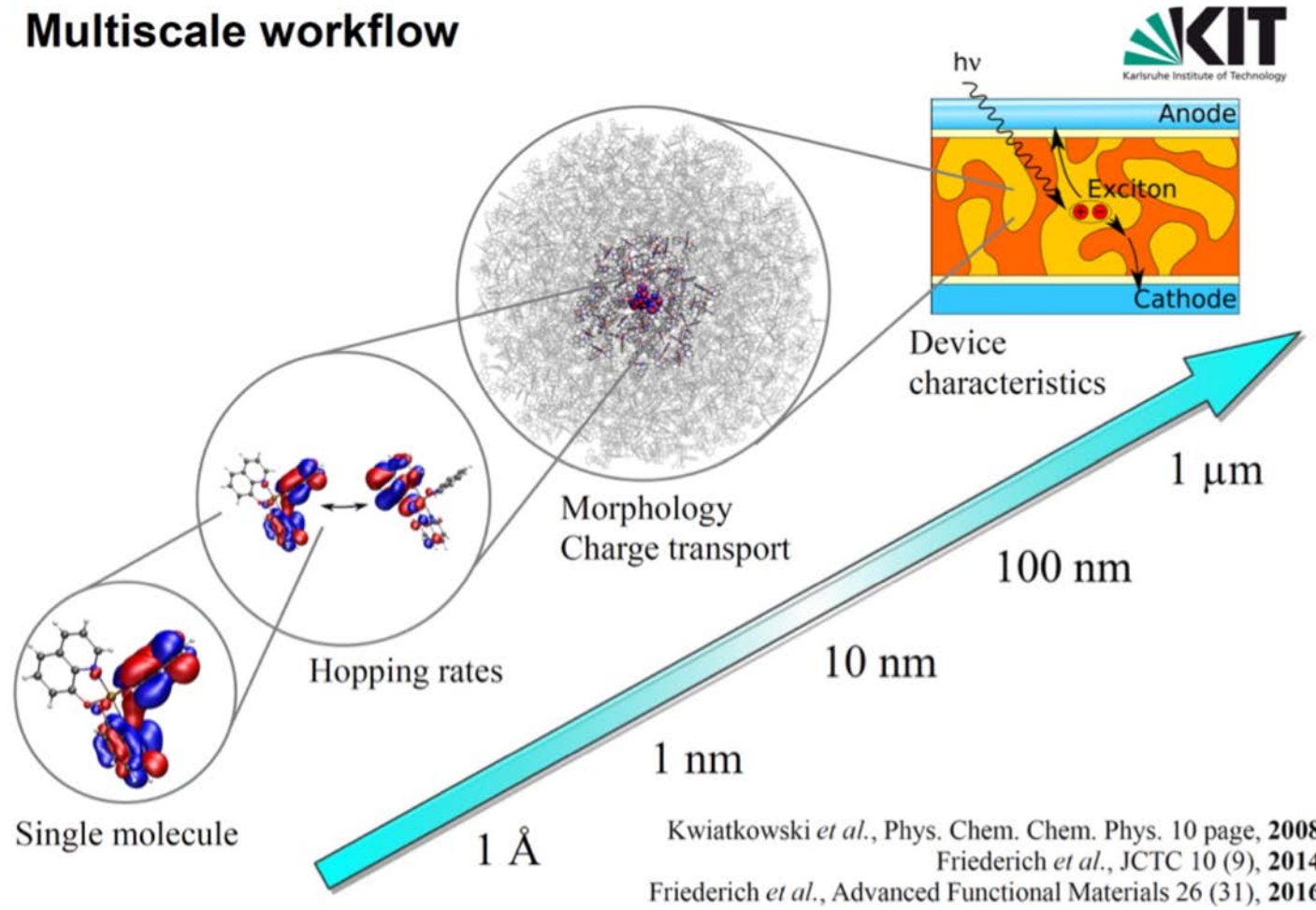
- Rapid Prototyping with Graphical Workflow editor (SimStack)



Prof. Dr.  
Wolfgang Wenzel

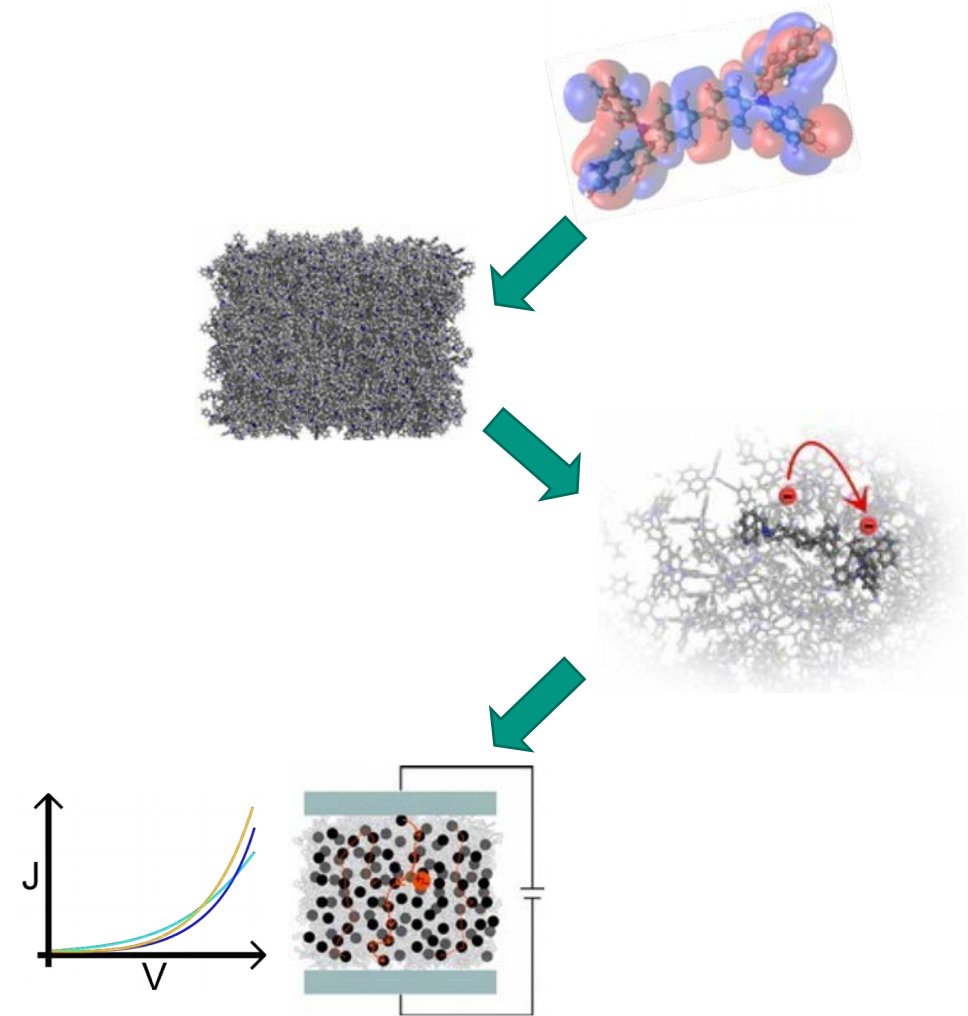
# Multiscale Materials Modelling and Virtual Design AG-Wenzel (INT)

## Multiscale workflow



# Multiscale simulation in Organic Electronics

- 1. Single molecule parametrization (QM)
  - Geometry optimization
  - Customized force-fields
- 2. Generation of atomistic morphologies
  - Molecules parametrized on quantum mechanical level
  - Simulation of physical vapor deposition
- 3. Calculation of charge hopping rates
  - Full quantum mechanical electronic structure analysis
  - Electronic couplings, reorganization and orbital energies
- 4. Charge transport simulations
  - Time resolved charge carrier/exciton dynamics
  - IVs, IQEs, carrier balance, quenching, ...



Content adapted from:



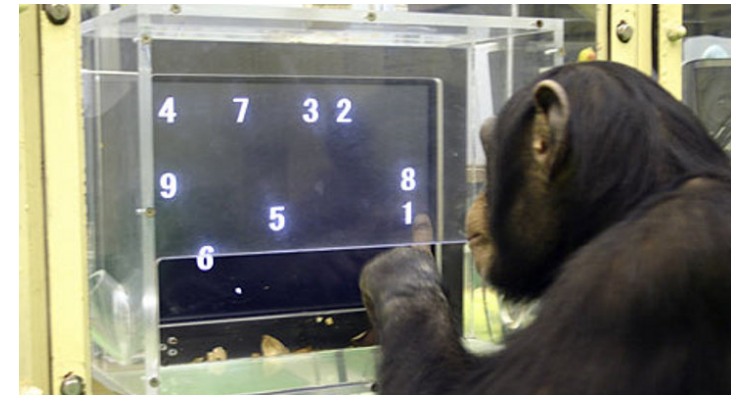
# Challenges and Opportunities of Workflow Technology

- **Multiscale Materials Modelling and Virtual Design @KIT**
- **Challenges in Materials Modelling**
- Opportunities of Workflow Technology
- Challenges of Workflow Technology

# Challenges - Integration of new employees

- New student in the office
  - Reproduction of a given result as first task
  - Application of the solution to another data set (Bachelor level)
  - Improvement of the given method (Master level)
  - Development of a new method (PhD Cand. level)
- Reality
  - Bachelor Thesis – 3 months
  - Barely enough time to get familiar with the work environment (command line, ssh, HPC, software)

Training of new Group Members

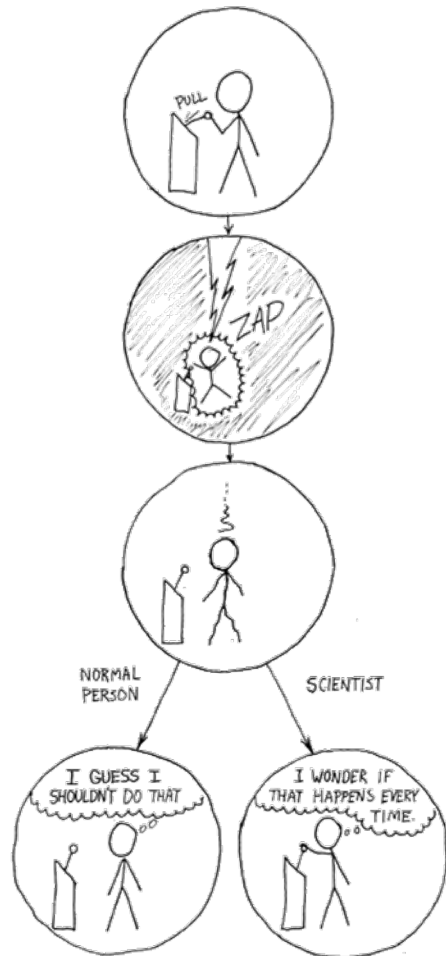


Content adapted from:





# Challenges – Reproducibility/Replicability



xkcd.com

## ■ Motivation

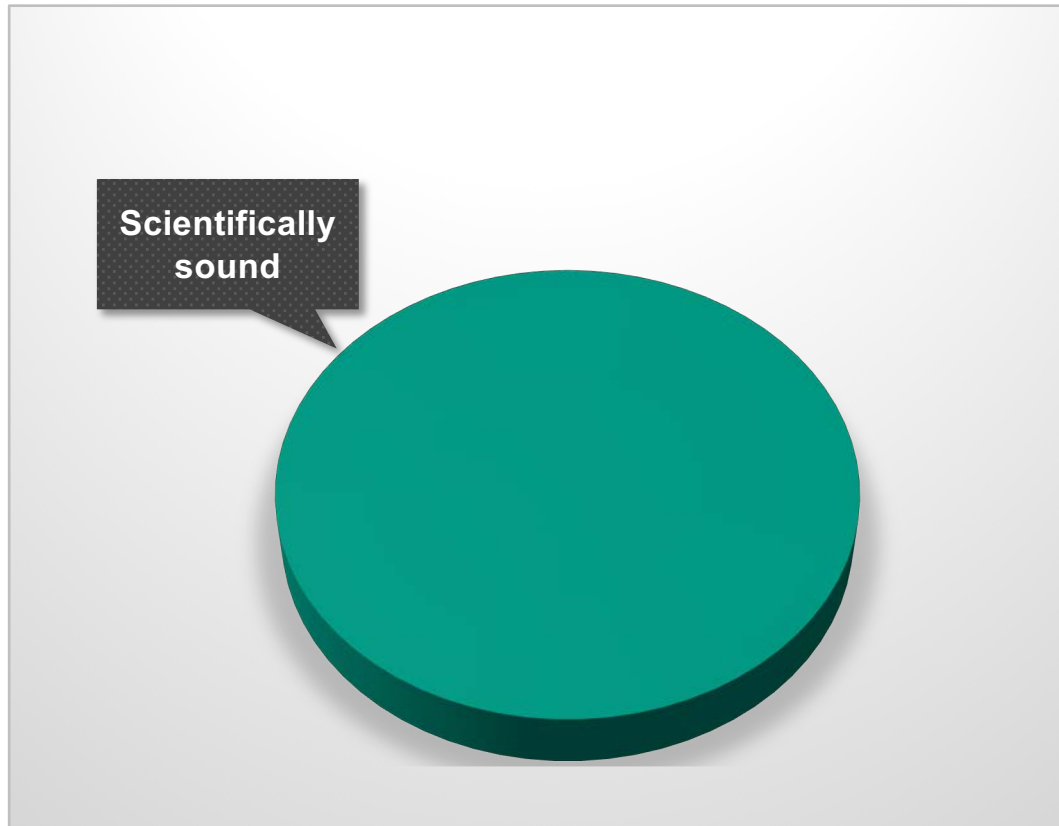
- Continue ongoing projects
- Build on work of others (Publications)
- Use existing data for Benchmarking/Datamining (Repository)
- Increase impact of own publications (citations)

## Reproducibility

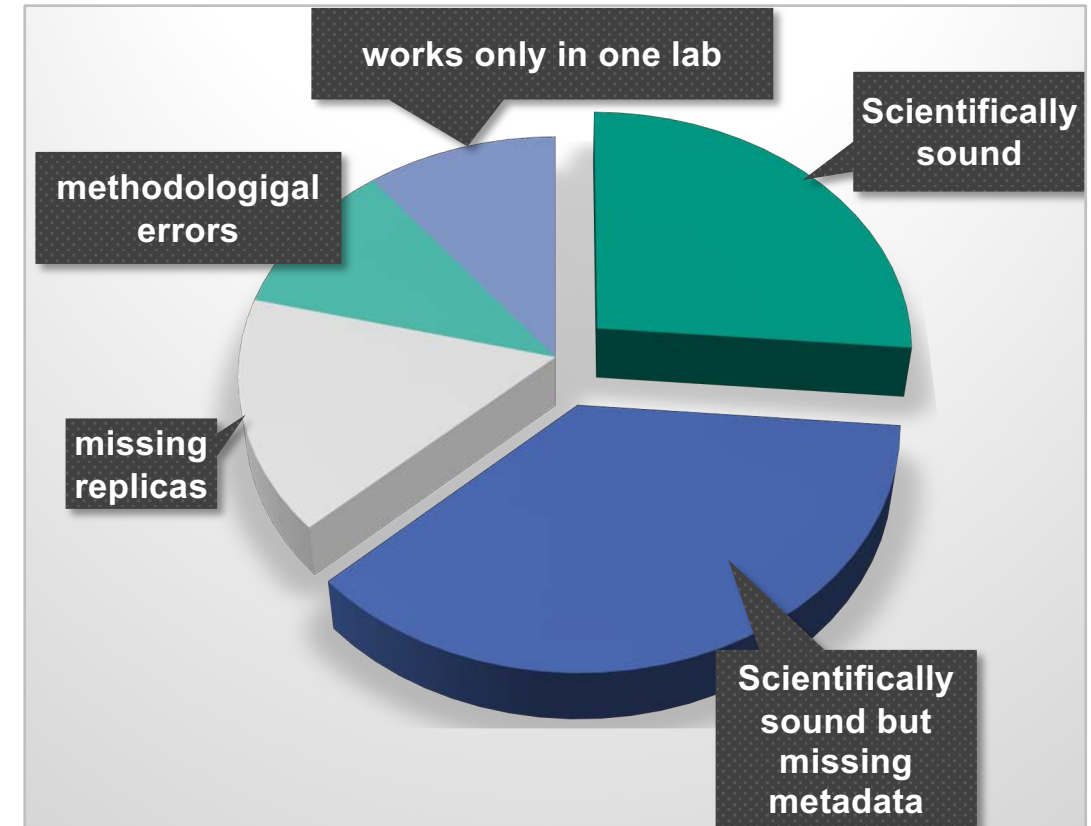


# Challenges – Reproducibility/Replicability

## Quality of Published Data (Journal/Database)



VS.



fictive data

# Challenges – Reproducibility/Replicability

## Molecular modeling

Homology models of the GPHR Chimeras ectodomains in complex with the hormones and their visual representations were generated using the Molecular Operating Environment (MOE, 2012.10; Chemical Computing Group Inc., Montreal, QC, Canada). For all homology models, the crystal structure of the FSHR ectodomain in complex with FSH (PDB code 4AY9) was used as template [8]. The protein sequences of the human TSHR and hTSH were acquired from the UniProt database (Accession number hTSHR: P16473, hTSH: P01222) [32] and chimera ectodomain sequences generated by combining the corresponding parts of TSHR and FSHR protein sequence. The sequences of the

chimeras were aligned to the crystal structure within MOE prior to simulation. For each chimera 100 homology models at 300 K were generated employing homology modeling the hormone was retained as environment. For the LRR domain whereas for the alpha domain, only the common environment while the beta domain the coordinates of the FSH

Schaarschmidt,  
PLOS ONE (2014)

## Experimental Section

*Simulations:* DFT calculations were carried using the Turbomole Package.<sup>[14]</sup> All calculations were performed using the hybrid B3-LYP<sup>[25]</sup> functional. Reorganization energies were calculated using the def2-TZVP<sup>[26]</sup> basis-set while for energy levels, energy disorder, and electronic couplings, the def2-SV(P)<sup>[27]</sup> basis-set was used. Atomistically resolved morphologies were generated using the Metropolis Monte Carlo based simulated annealing method DEPOSIT.<sup>[13]</sup> This method required

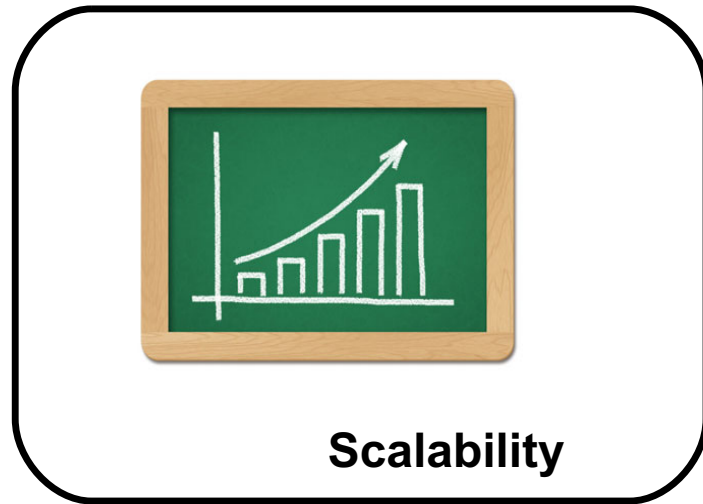
DFT-optimized molecular conformations and partial charges (B3-LYP/def2-SV(P)). Energy disorder and HOMO/LUMO levels as well as IPs and EA were calculated using the Quantum Patch method.<sup>[16]</sup>

Friederich, *Advanced Materials* (2017)

- Preparation of input data
- Missing steps/ customizations
  - e.g. Data sources
- Hardware
- Compiler options
- Software Versions
- ...

→ For Reproducibility, all important aspects of a simulation need to be captured

# 4. Challenges - Scalability



- Requirements
  - Availability of Resources
  - Automatization of the simulation protocol
  - Homogeneous input data

HOW LONG CAN YOU WORK ON MAKING A ROUTINE TASK MORE EFFICIENT BEFORE YOU'RE SPENDING MORE TIME THAN YOU SAVE?  
(ACROSS FIVE YEARS)

HOW OFTEN YOU DO THE TASK

|            | 50/DAY   | 5/DAY     | DAILY      | WEEKLY     | MONTHLY    | YEARLY     |
|------------|----------|-----------|------------|------------|------------|------------|
| 1 SECOND   | 1 DAY    | 2 HOURS   | 30 MINUTES | 4 MINUTES  | 1 MINUTE   | 5 SECONDS  |
| 5 SECONDS  | 5 DAYS   | 12 HOURS  | 2 HOURS    | 21 MINUTES | 5 MINUTES  | 25 SECONDS |
| 30 SECONDS | 4 WEEKS  | 3 DAYS    | 12 HOURS   | 2 HOURS    | 30 MINUTES | 2 MINUTES  |
| 1 MINUTE   | 8 WEEKS  | 6 DAYS    | 1 DAY      | 4 HOURS    | 1 HOUR     | 5 MINUTES  |
| 5 MINUTES  | 9 MONTHS | 4 WEEKS   | 6 DAYS     | 21 HOURS   | 5 HOURS    | 25 MINUTES |
| 30 MINUTES |          | 6 MONTHS  | 5 WEEKS    | 5 DAYS     | 1 DAY      | 2 HOURS    |
| 1 HOUR     |          | 10 MONTHS | 2 MONTHS   | 10 DAYS    | 2 DAYS     | 5 HOURS    |
| 6 HOURS    |          |           |            | 2 MONTHS   | 2 WEEKS    | 1 DAY      |
| 1 DAY      |          |           |            |            | 8 WEEKS    | 5 DAYS     |

HOW MUCH TIME YOU SHAVE OFF

xkcd.com

# Challenges - Competence Drain

- Highly specialized employee (PhD cand., etc.) implements method
  - Results in highly specialized hard to use software tool
- Employee leaves
  - Knowledge about usage of software tool leaves with employee
  - Usage/Maintenance/Support of software tool hard to impossible



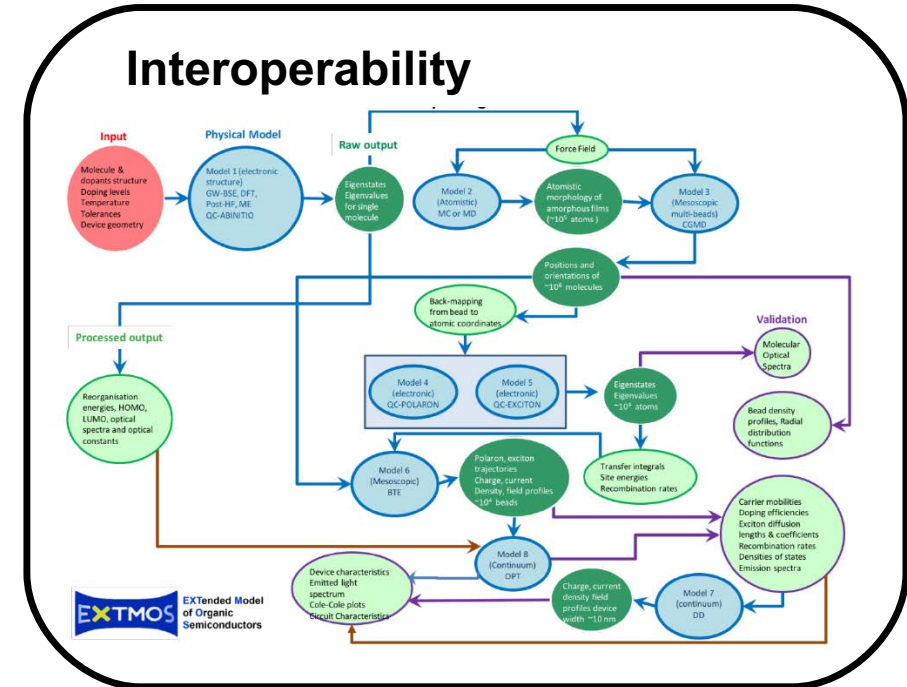
*//Peter wrote this, nobody knows what it does, don't change it!*

Content adapted from:



# Challenges - Interoperability

- Software needs to be inter-operable
  - Multi-scale environments
  - Data-sharing
  
- Software Development
  - Software needs to be usable by other groups
  - Software aggregates need to be
    - prepared
    - shared
    - archived
    - presented (Deliverables)



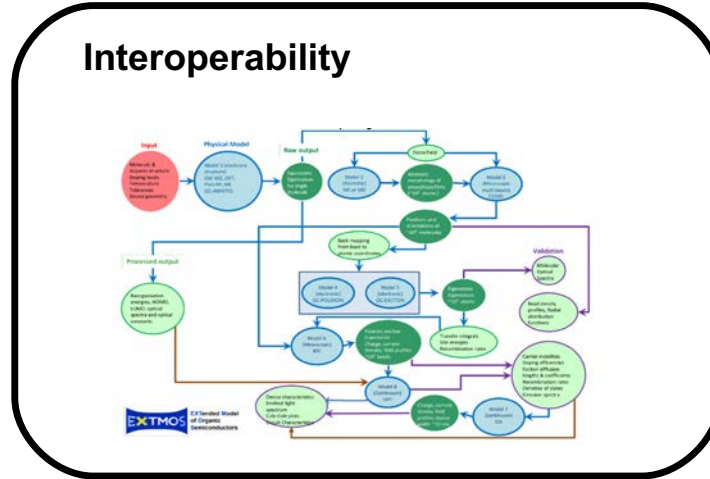
*//When I wrote this, only God and I understood what I was doing  
 //Now, God only knows*

Content adapted from:




# Challenges – Scientific Group

**Competence Drain**


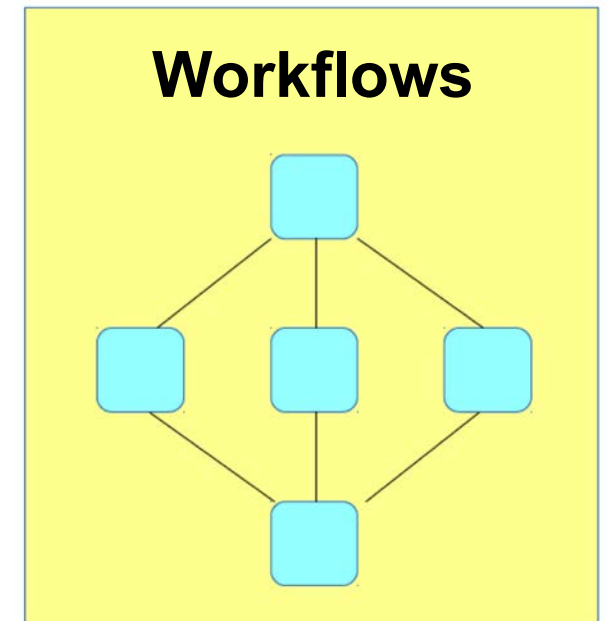



**Reproducibility**



**Scalability**

**Training of new group members**

Content adapted from: 

# Challenges and Opportunities of Workflow Technology

- **Multiscale Materials Modelling and Virtual Design @KIT**
- **Challenges in Materials Modelling**
- **Opportunities of Workflow Technology**
- **Challenges of Workflow Technology**



# Workflows

| Focus  | Components   | Abstraction Level   | Backend   |
|--|--|---|---|
| <ul style="list-style-type: none"><li>• Automation</li><li>• Speedup</li><li>• Provenance</li><li>• Innovation</li></ul> | <ul style="list-style-type: none"><li>• Applications</li><li>• (Web)services/Utilities</li><li>• Libraries</li><li>• Data</li><li>• Control Elements</li></ul> | <ul style="list-style-type: none"><li>• Command Line</li><li>• Script based</li><li>• Workflow file</li><li>• GUI</li></ul> | <ul style="list-style-type: none"><li>• HPC Clusters</li><li>• Cloud/Grid Resources</li><li>• Workstation</li></ul> |

<https://github.com/MD-Studio/MDStudio/>

A workflow represents the coordinated execution of repeatable computational steps while accounting for dependencies and concurrency of tasks.

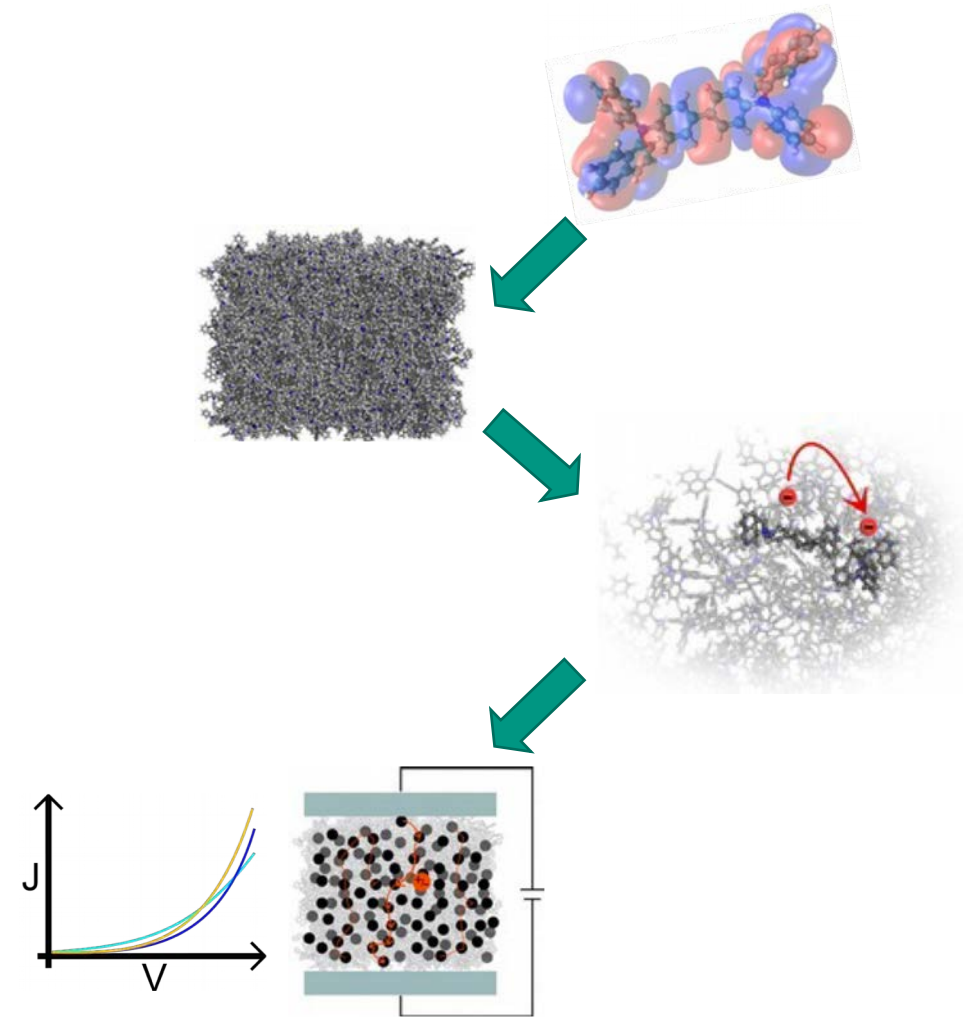
# Workflows

| Focus  | Components   | Abstraction Level   | Backend   |
|--|--|---|---|
| <ul style="list-style-type: none"><li>• Automation</li><li>• Speedup</li><li>• Provenance</li><li>• Innovation</li></ul> | <ul style="list-style-type: none"><li>• Applications</li><li>• (Web)services/Utilities</li><li>• Libraries</li><li>• Data</li><li>• Control Elements</li></ul> | <ul style="list-style-type: none"><li>• Command Line</li><li>• Script based</li><li>• Workflow file</li><li>• GUI</li></ul> | <ul style="list-style-type: none"><li>• HPC Clusters</li><li>• Cloud/Grid Resources</li><li>• Workstation</li></ul> |



# The multiscale simulation workflow for Organic Electronics

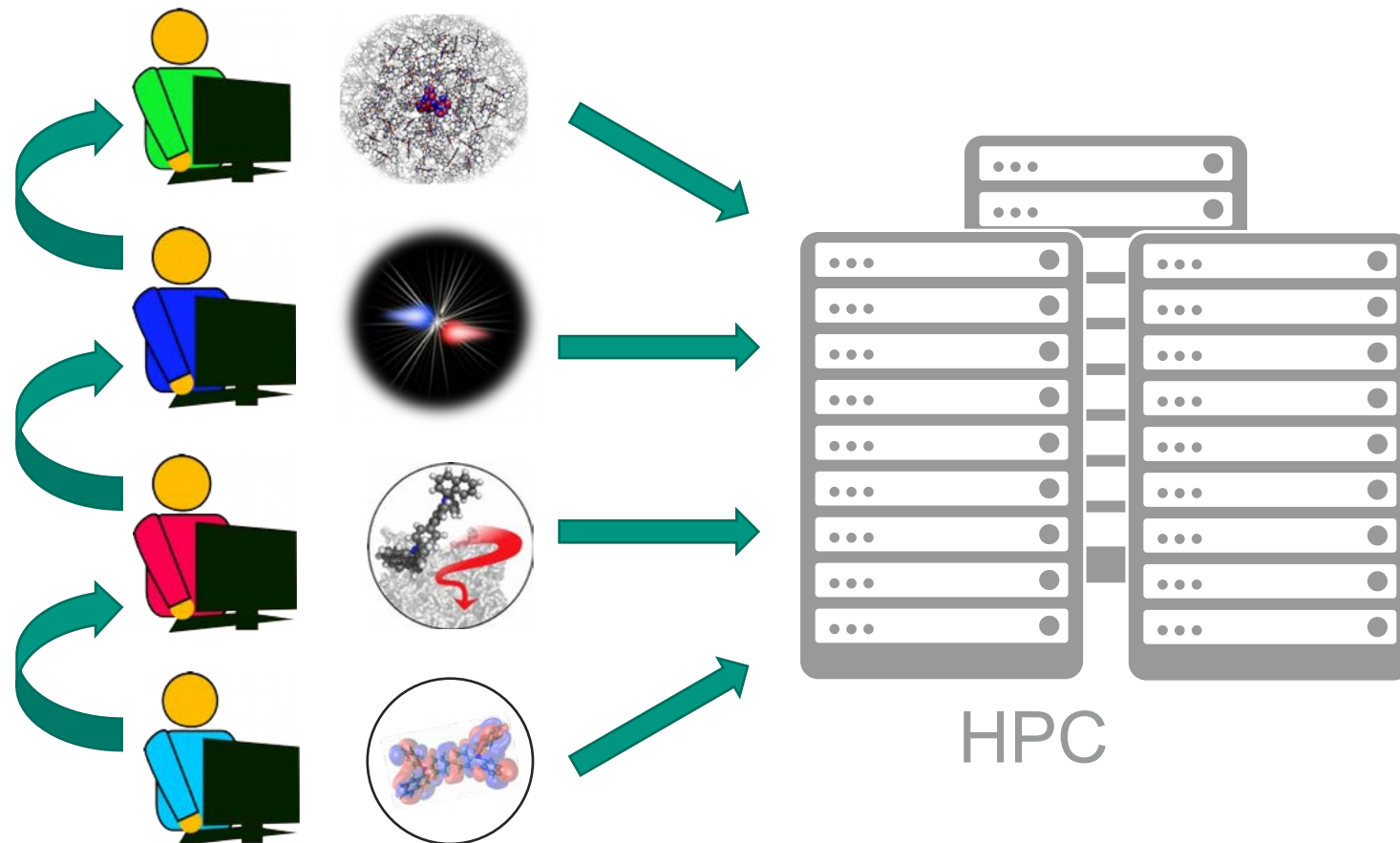
- 1. Single molecule parametrization (QM)
  - Geometry optimization
  - Customized force-fields
- 2. Generation of atomistic morphologies
  - Molecules parametrized on quantum mechanical level
  - Simulation of physical vapor deposition
- 3. Calculation of charge hopping rates
  - Full quantum mechanical electronic structure analysis
  - Electronic couplings, reorganization and orbital energies
- 4. Charge transport simulations
  - Time resolved charge carrier/exciton dynamics
  - IVs, IQEs, carrier balance, quenching, ...



Content adapted from:



# The multiscale simulation workflow for Organic Electronics



- High level of complexity of multiscale modeling:
  - Input preparation
  - Data transfer to computing resources
  - Job submission and monitoring,

→ Traditional approach to multiscale modeling

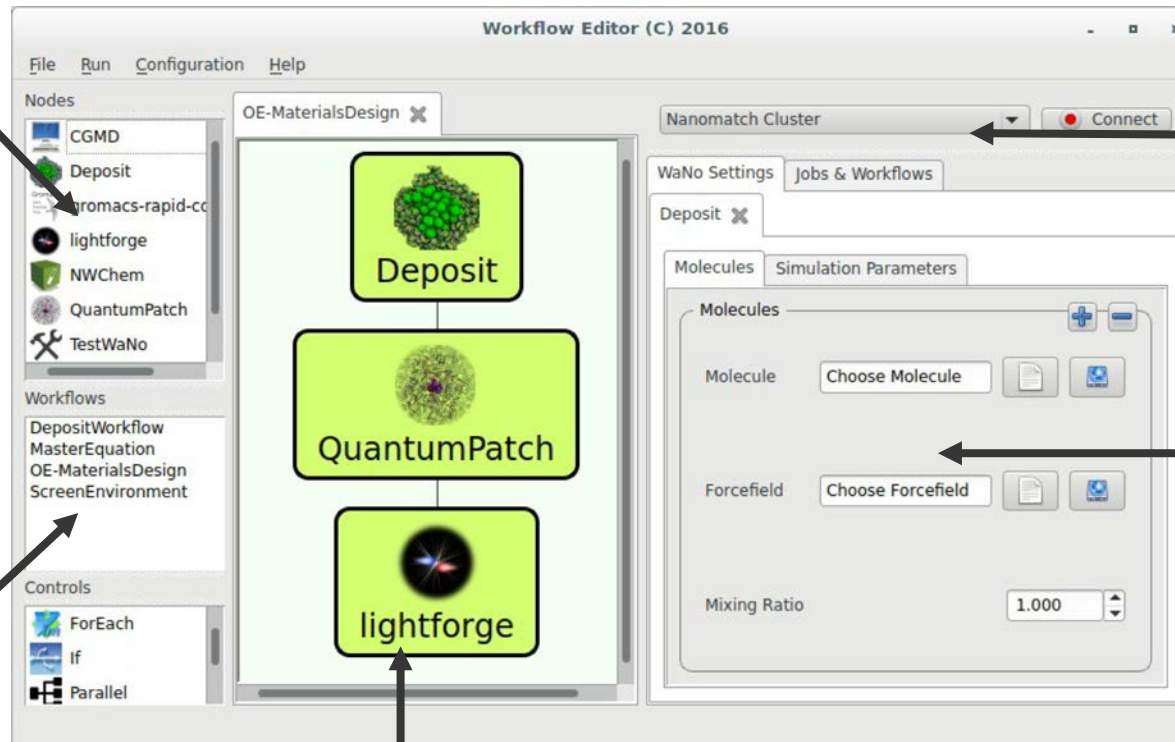
Content adapted from:



# Translation enabled by SimStack



Embedded Scientific modules = „WaNos“



Connect to remote computational resources

Define input files and parameters for each module

Saved workflows for reproducible multiscale simulations

Construct a workflow by drag & drop



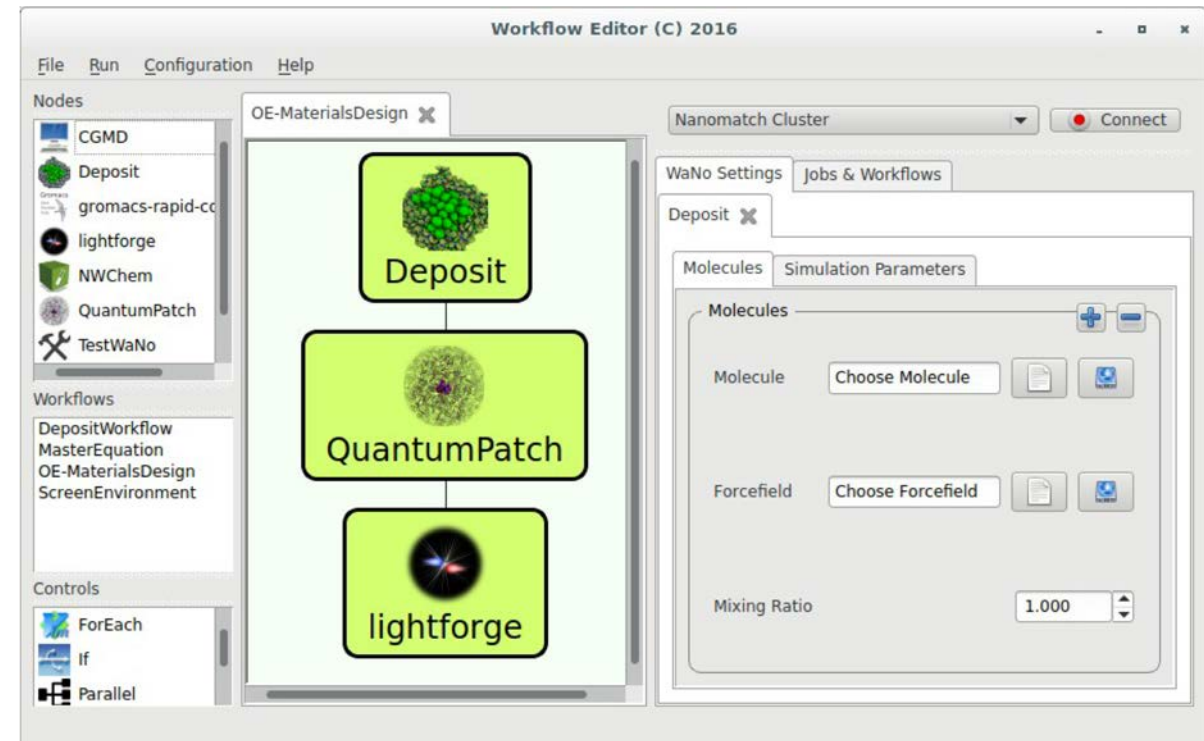
Content adapted from:

# Translation enabled by SimStack

... a generic workflow platform conquering complexity



- Open to arbitrary software modules
- Rapid prototyping: 30min to include new modules, 1 h to construct functional workflows
- Maximal reusability and scalability
- Module interoperability:
  - Fully automated, file based data transfer between modules
  - Schema based data transfer in development
  - Compatible with OWL ontologies (e.g. EMMO, once developed)



Content adapted from:



# Workflows @INT - Example Deposit

## Deposit

- Deposit is a Monte-Carlo tool to generate organic thin-film morphologies
  - Complex input language (roughly 80 parameters up front)
  - Minimum number of input files: 2
  - Minimum number of parameters, which you have to actually know: 7

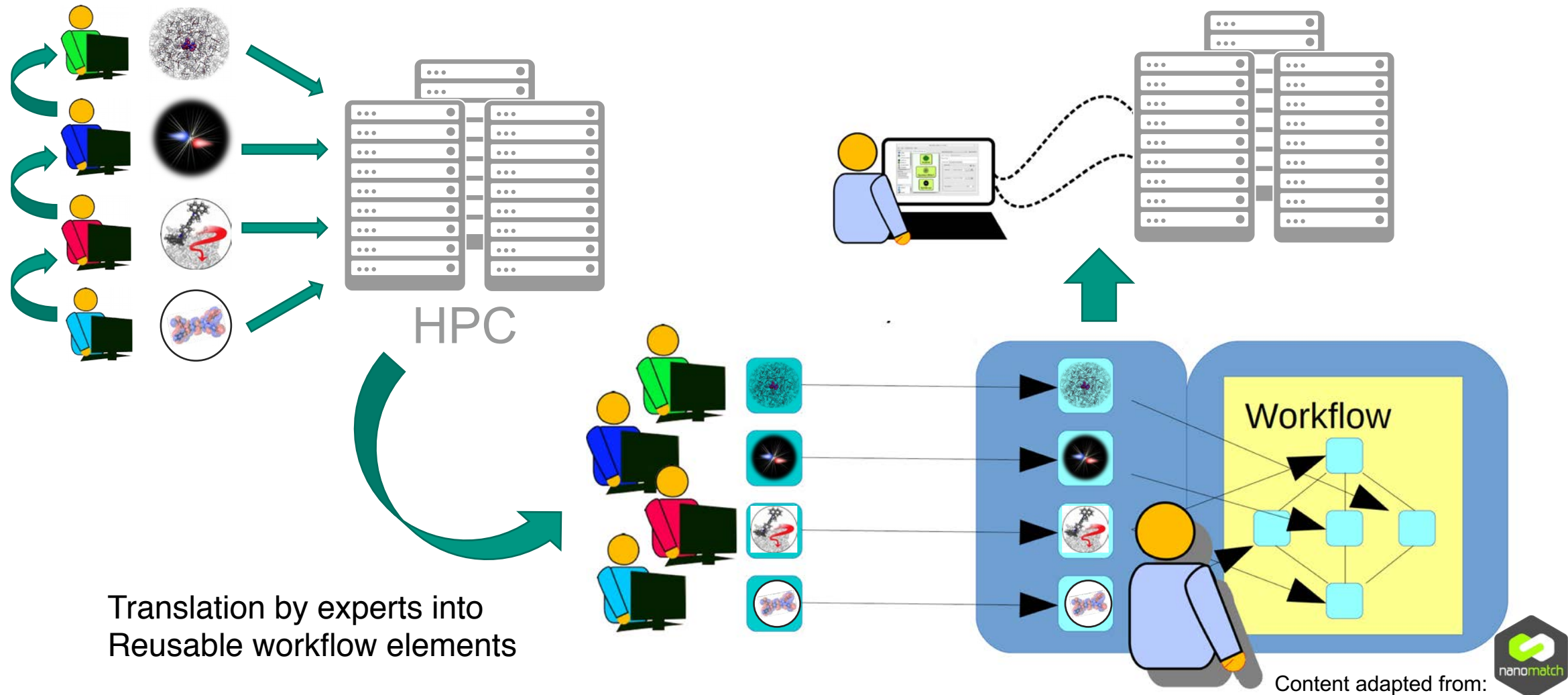
## Learning curve

- Learn bash
- Find out about the important parameters
- Learn ssh
- Scp files over
- Learn specific qsub commands
- Write (or at least adapt) a submission script
- ...

Content adapted from:



# The multiscale simulation workflow for Organic Electronics





# Incorporation of Modules (Workflow Active Nodes - WaNos)

```
<WaNoTemplate>  
<WaNoRoot name="Your Wano Name Here" >  
  
  Input parameters here  
  
</WaNoRoot>  
<WaNoExecCommand>./ExecutionScript.sh</WaNoExecCommand>  
  Command for program execution  
  
<WaNoInputFiles>  
  <WaNoInputFile logical_filename="ExecutionScript.sh">ExecutionScript.sh</WaNoInputFile>  
  <WaNoInputFile logical_filename="MandatoryInputFile.dat">MandatoryInputFile</WaNoInputFile>  
  Mandatory input files  
</WaNoInputFiles>  
<WaNoOutputFiles>  
  <WaNoOutputFile>MandatoryOutputFileName1.dat</WaNoOutputFile>  
  <WaNoOutputFile>MandatoryOutputFileName2.dat</WaNoOutputFile>  
  File output expected  
  by other WaNos  
</WaNoOutputFiles>  
</WaNoTemplate>
```

Content adapted from:



# Incorporation of Modules (Workflow Active Nodes - WaNos)

```

<WaNoTemplate>
<WaNoRoot name="Random Number Generator">

  <WaNoFloat name="This is a float field">0.0</WaNoFloat>

  <WaNoInt name="This is an int field">1</WaNoInt>

  <WaNoChoice name="Choice of options">
    <Entry id="0" chosen="True">Standard option</Entry>
    <Entry id="1">Second option</Entry>
    <Entry id="2">Third option</Entry>
  </WaNoChoice>

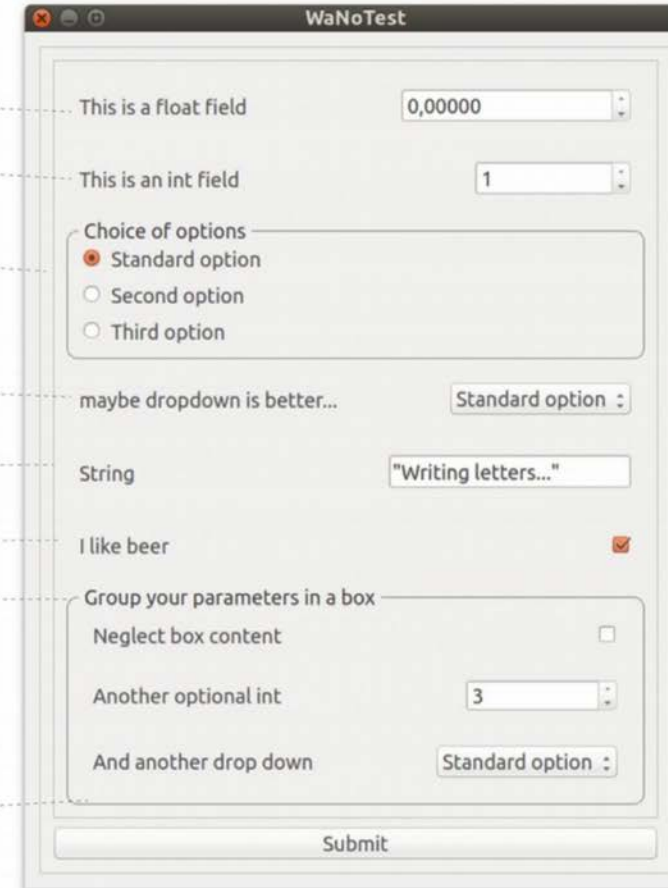
  <WaNoDropDown name="maybe dropdown is better...">
    <Entry id="0" chosen="True">Standard option</Entry>
    <Entry id="1">Second option</Entry>
    <Entry id="2">Third option</Entry>
  </WaNoDropDown>

  <WaNoString name="String">"Writing letters..."</WaNoString>

  <WaNoBool name="I like beer">True</WaNoBool>

  <WaNoBox name="Group your parameters in a box">
    <WaNoBool name="Neglect box content">False</WaNoBool>
    <WaNoInt name="Another optional int">3</WaNoInt>
    <WaNoDropDown name="And another drop down">
      <Entry id="0" chosen="True">Standard option</Entry>
      <Entry id="1">Second option</Entry>
      <Entry id="2">Third option</Entry>
    </WaNoDropDown>
  </WaNoBox>

```



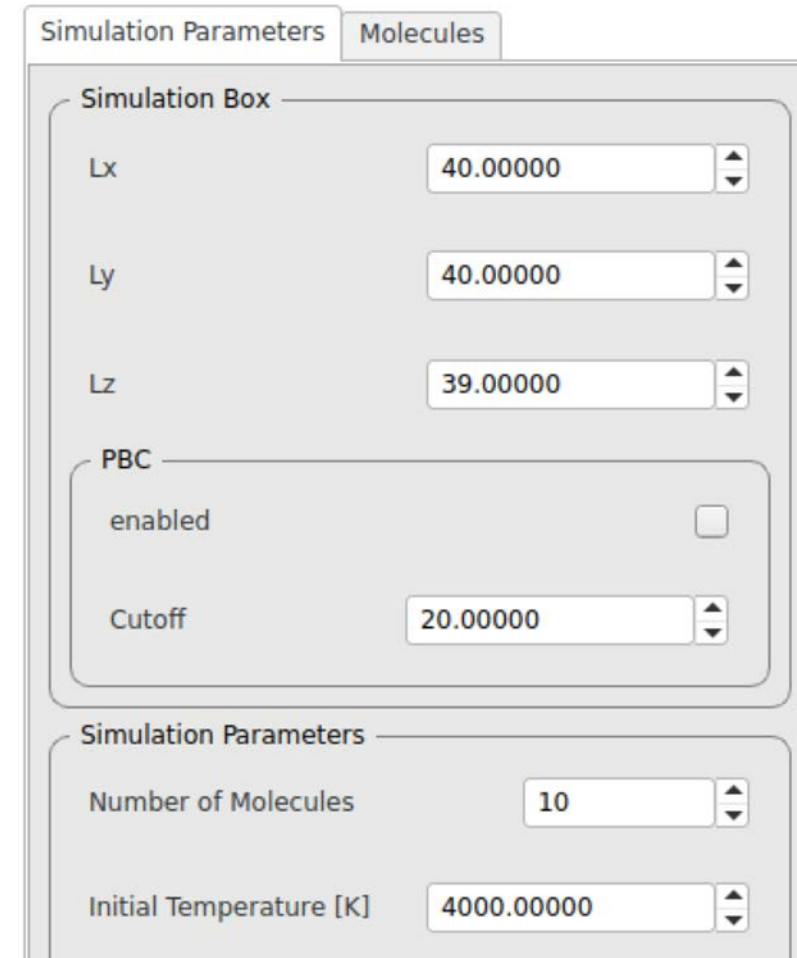
The screenshot shows a window titled "WaNoTest" containing a form with the following elements:

- "This is a float field" with a text input containing "0,00000".
- "This is an int field" with a text input containing "1".
- "Choice of options" with three radio buttons: "Standard option" (selected), "Second option", and "Third option".
- "maybe dropdown is better..." with a dropdown menu showing "Standard option".
- "String" with a text input containing "Writing letters...".
- "I like beer" with a checked checkbox.
- "Group your parameters in a box" containing:
  - "Neglect box content" with an unchecked checkbox.
  - "Another optional int" with a text input containing "3".
  - "And another drop down" with a dropdown menu showing "Standard option".
- A "Submit" button at the bottom.

# Workflows @INT - Example Deposit

## Deposit

- Deposit is a Monte-Carlo tool to generate organic thin-film morphologies
  - Complex input language (roughly 80 parameters up front)
  - Minimum number of input files: 2
  - Minimum number of parameters, which you have to actually know: 7



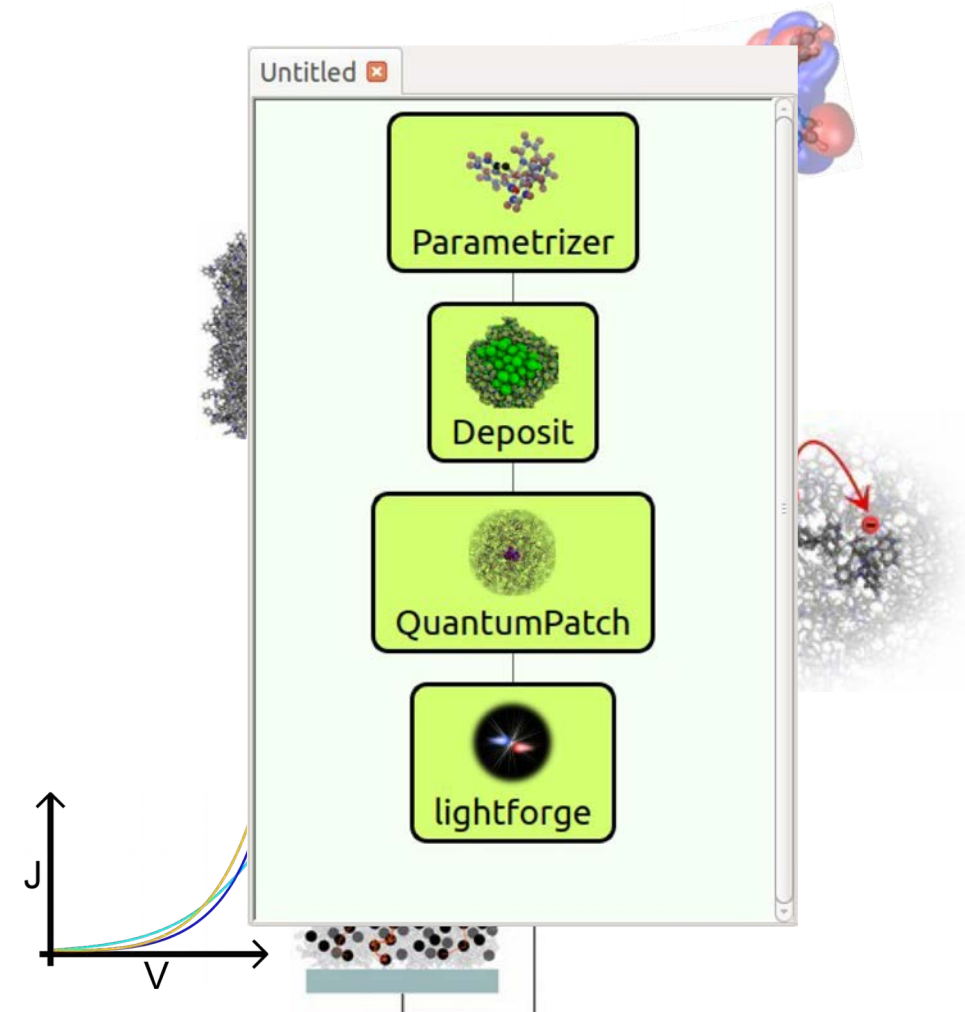
The screenshot shows the 'Simulation Parameters' and 'Molecules' tabs. The 'Simulation Box' section includes Lx (40.00000), Ly (40.00000), and Lz (39.00000). The 'PBC' section has an 'enabled' checkbox (unchecked) and a 'Cutoff' value of 20.00000. The 'Simulation Parameters' section includes 'Number of Molecules' (10) and 'Initial Temperature [K]' (4000.00000).

Content adapted from:



# The multiscale simulation workflow for Organic Electronics

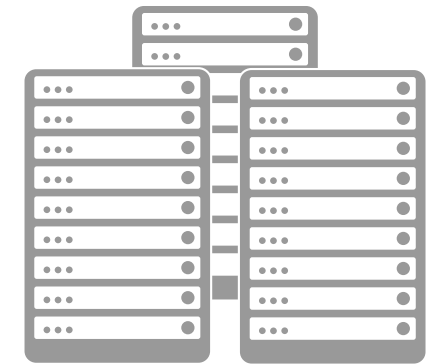
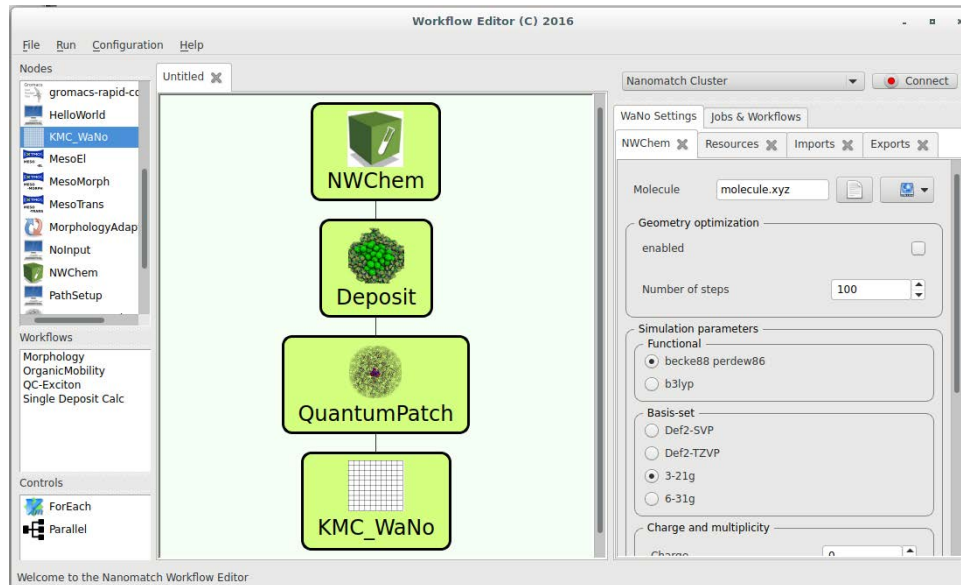
- 1. Single molecule parametrization (QM)
  - Geometry optimization
  - Customized force-fields
- 2. Generation of atomistic morphologies
  - Molecules parametrized on quantum mechanical level
  - Simulation of physical vapor deposition
- 3. Calculation of charge hopping rates
  - Full quantum mechanical electronic structure analysis
  - Electronic couplings, reorganization and orbital energies
- 4. Charge transport simulations
  - Time resolved charge carrier/exciton dynamics
  - IVs, IQEs, carrier balance, quenching, ...



Content adapted from:



# Connection to Compute Resources



HPC

- Middleware actively developed in FZ Jülich
- Handles User Authentication
- Can connect to all common schedulers
- Handles data transfer between workflow steps
- Setup in < 1h with installer provided by Nanomatch

Content adapted from:



# Challenges and Opportunities of Workflow Technology

- **Multiscale Materials Modelling and Virtual Design @KIT**
- **Challenges in Materials Modelling**
- **Opportunities of Workflow Technology**
- **Challenges of Workflow Technology**

# Reproducibility of Workflows

What needs to be captured?

- Software, Tools, Scripts
- WaNos
- Workflow templates
- Executed Workflows
- Author



How to capture it

- Containerization (e.g. udocker)



<https://github.com/indigo-dc/udocker>

- Versioning 

- Tags 

- Repositories 

# Interoperability of workflow frameworks



It is unlikely that one workflow framework will meet all the needs of the community

→ Shared data Formats and/or Converters will be required



# Acknowledgements

- Prof. Dr. Wolfgang Wenzel
- Dr. Ivan Kondov (SCC)



- Dr. Timo Strunk
- Dr. Tobias Neumann

<http://www.nanomatch.com>

<http://www.simstack.eu>

# Thank You!



joerg.schaarschmidt@kit.edu

