

<sup>1</sup>

## Materials Genomics

<sup>2</sup>

**Philipp Pelz**

---

Corresponding author: Philipp Pelz,

3   **Abstract**

4   This course introduces students to materials genomics, treating the periodic table  
5   and the space of known crystal structures as a searchable, computable design space.  
6   Students learn how materials databases are built, how simulation methods generate  
7   materials data, how atomic structure is represented numerically, how structure–  
8   property relationships are learned using machine learning, and how uncertainty-  
9   aware models enable accelerated materials discovery.

10   **1 Course Information**

11   **4th/5th Semester – 5 ECTS · 2h lecture + 2h exercises per week**

12   *Coordinated with “Mathematical Foundations of AI & ML” (MFML) and*  
13   *“ML for Materials Processing & Characterization” (ML-PC)*

---

14   **2 Course Philosophy**

15   Materials genomics views the periodic table and all known crystal structures as a  
16   **high-dimensional design space.**

17   In this course, students learn to:

- 18   • understand how materials data is generated by simulations and experiments,  
19   • treat materials data as a structured, learnable representation space,  
20   • move beyond classical descriptors toward learned representations,  
21   • use ML models as surrogates for quantum-mechanical and continuum simula-  
22   tions,  
23   • reason about uncertainty, stability, and discovery,  
24   • understand how computational screening integrates with experiments.

25   The course explicitly **builds on MFML**:

- 26   • PCA and regression are assumed background,  
27   • neural networks, representation learning, and uncertainty are used, not re-  
28   derived.
- 

29   **3 Week-by-Week Curriculum (14 weeks)**

30   **3.1 Unit I — Where Materials Data Comes From (Weeks 1–4)**

31   **3.1.1 Week 1 – What is Materials Genomics?**

- 32   • Genomics analogy: genes → functions vs atoms → properties.  
33   • Structure–property–processing paradigm from a *structure-first* viewpoint.  
34   • Materials databases as design spaces: Materials Project, OQMD, AFLOW,  
35   NOMAD.

36   **Exercise:**

37   Explore Materials Project; query bandgaps, formation energies, symmetries.

---

38   **3.1.2 Week 2 – Simulation methods as data generators**

- 39   • Why simulations dominate materials data generation.  
40   • Simulation methods as mappings from assumptions to data.  
41   • Overview of scales and outputs:  
42    – FEM: continuum fields (stress, strain).  
43    – MD: trajectories, forces, diffusion.  
44    – MC: thermodynamic sampling.  
45    – DFT: energies, electronic structure.  
46   • Accuracy–cost–scale trade-offs and systematic biases.

**Exercise:**

For selected materials properties, identify suitable simulation methods and expected biases.

---

**3.1.3 Week 3 – Atomistic and electronic simulations (DFT, MD, MC)**

- Density Functional Theory: ground-state bias, exchange–correlation functionals, consistency vs accuracy.
- Molecular Dynamics: force fields, time averaging, limitations of timescales.
- Monte Carlo: phase-space sampling and thermodynamic averages.
- What quantities in materials databases come directly from simulations.

**Exercise:**

Inspect Materials Project entries; identify simulation assumptions and derived quantities.

---

**3.1.4 Week 4 – Continuum simulations, thermodynamics, and stability**

- FEM as a structure–property mapping at the continuum scale.
- Constitutive models as implicit surrogates.
- Formation energies, convex hulls, metastability.
- Why “stable” does not imply “synthesizable”.

**Exercise:**

Analyze stability and simulated properties for a small materials system.

---

**3.2 Unit II — Representations of Materials (Weeks 5–7)**

*(Aligned with early neural networks in MFML)*

**3.2.1 Week 5 – From classical descriptors to learned representations**

- Classical descriptors: Magpie, matminer.
- Limits of hand-crafted features.
- Motivation for representation learning.

**Exercise:**

Build a simple property predictor using classical descriptors.

---

**3.2.2 Week 6 – Graph-based crystal representations**

- Crystals as graphs: nodes, edges, periodic boundary conditions.
- Intuition behind CGCNN and MEGNet.
- Relation to neural network concepts from MFML.

**Exercise:**

Construct and visualize graph representations of crystal structures.

---

**3.2.3 Week 7 – Local atomic environments**

- Local vs global representations.
- Coordination environments, Voronoi tessellations.
- SOAP descriptors as a bridge to learned representations.

**Exercise:**

Compute SOAP vectors and explore similarity in descriptor space.

---

95    **3.3 Unit III — Learning Structure–Property Relations (Weeks 8–10)**

96    **3.3.1 Week 8 – *Regression and generalization in materials data***

- 97    • Predicting bandgaps, elastic moduli, formation energies.  
98    • Bias–variance trade-off and overfitting.  
99    • Dataset size vs model complexity.

100    **Exercise:**

101    Compare linear, random forest, and neural network regressors.

---

102    **3.3.2 Week 9 – *Neural networks for materials properties***

- 103    • Neural networks as surrogates for DFT-level properties.  
104    • Training pitfalls: data leakage, imbalance, extrapolation.  
105    • Interpretability challenges.

106    **Exercise:**

107    Train a small neural network and analyze generalization behavior.

---

109    **3.3.3 Week 10 – *Representation learning and feature discovery***

- 110    • Learned vs engineered features.  
111    • Transferability across chemical systems.  
112    • What networks learn about chemistry and structure.

113    **Exercise:**

114    Compare model performance using raw descriptors vs learned embeddings.

---

116    **3.4 Unit IV — Latent Spaces, Uncertainty, and Discovery (Weeks 11–13)**

117    **3.4.1 Week 11 – *Latent spaces of materials***

- 118    • Autoencoders and embeddings for crystal data.  
119    • Interpreting latent dimensions.  
120    • Structure families and chemical intuition.

121    **Exercise:**

122    Train an autoencoder; visualize latent materials space.

---

124    **3.4.2 Week 12 – *Clustering, uncertainty, and discovery logic***

- 125    • Why clustering is not discovery.  
126    • Outliers, anomalies, and candidate identification.  
127    • Aleatoric vs epistemic uncertainty.

128    **Exercise:**

129    Contrast clustering results with latent-space exploration.

---

131    **3.4.3 Week 13 – *Uncertainty-aware discovery and Gaussian Processes***

- 132    • Gaussian Process regression as a gold standard for uncertainty.  
133    • Exploration vs exploitation.  
134    • Relevance to materials acceleration platforms.

135    **Exercise:**

136    Compare GP regression and neural network ensembles for screening tasks.

---

139      **3.5 Unit V — Constraints, Trust, and Synthesis (Week 14)**

140      **3.5.1 Week 14 – *Physical constraints, limits, and outlook***

- 141      • Stability, charge neutrality, and symmetry constraints.  
142      • Physics-informed learning in materials discovery.  
143      • What ML can and cannot discover.  
144      • Integration with experimental workflows.

145      **Exercise:**

146      Mini-project synthesis and presentation.

---

148      **4 Learning Outcomes**

149      Students completing this course will be able to:

- 150      • Explain how simulation methods generate materials data and introduce bias.  
151      • Navigate and interrogate major materials databases.  
152      • Represent crystal structures using descriptors, graphs, and learned embed-  
153      dings.  
154      • Train and evaluate ML models for predicting materials properties.  
155      • Understand latent spaces and their role in materials discovery.  
156      • Quantify and interpret uncertainty in materials predictions.  
157      • Apply ML responsibly to accelerate materials screening.  
158      • Critically assess the limits of data-driven materials discovery.