

Large Synthetic Datasets for Machine Learning Applications in Power Transmission Grids

Marc Gillioz



April 15, 2025

Working Group – Foundation Models for the Electric Grid

The Team



Institute of Energy and Environment
**University of Applied Sciences and
Arts of Western Switzerland (Sion)**

- **Prof. Philippe Jacquod**
Head of Energy
Efficiency Group
- **Dr. Marc Gillioz**
Senior scientist
- **Guillaume Dubuis**
Master's student



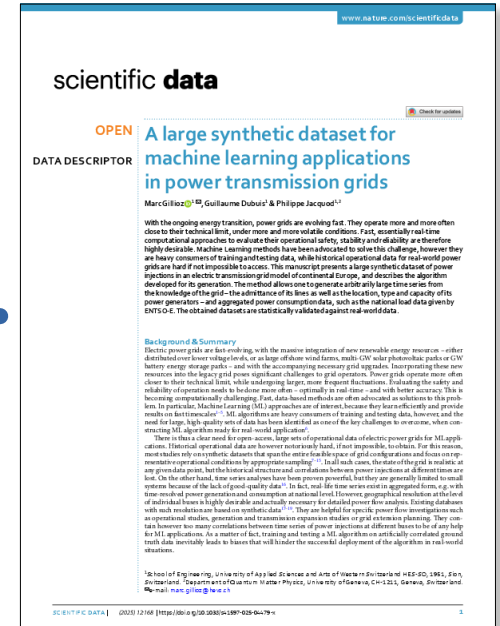
Motivation

- Power grids are under pressure → modeling needed
- Why transmission grids?
 - meshed and complex power flow, with universal features → interesting
 - limited losses at ultra-high voltage → DC power flow approx. valid
- armasuisse Cyber-Defence Campus mandate
 - detection of false data injection attacks



Outline

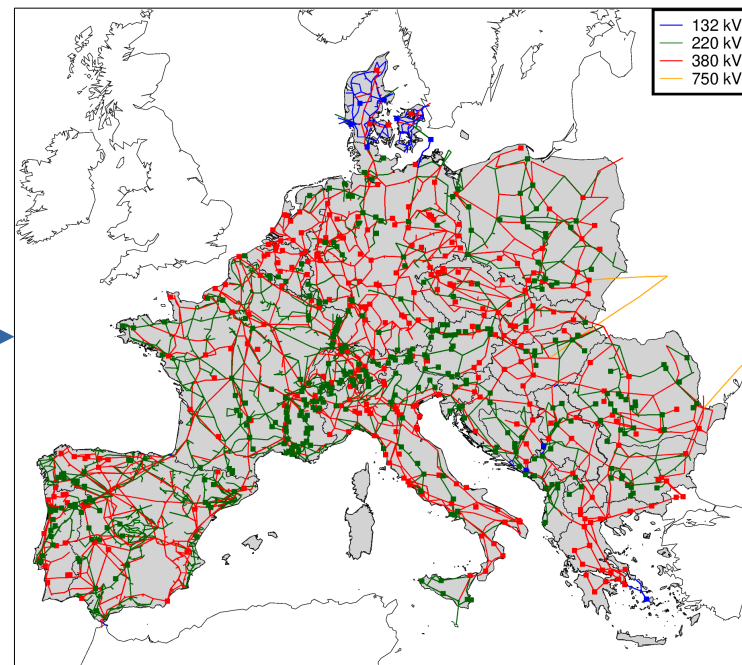
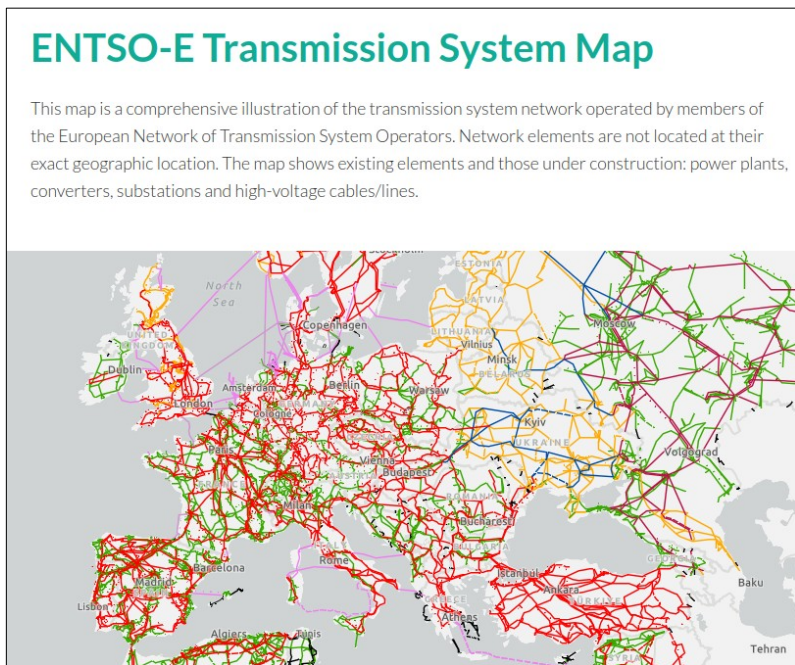
- The problem realistic time series for large grid models
- Our solution load modeling and optimal dispatch
- An application ML algorithms for anomaly detection



coming soon...

The problem: transmission grid models

- Limited access to real-world data. But...



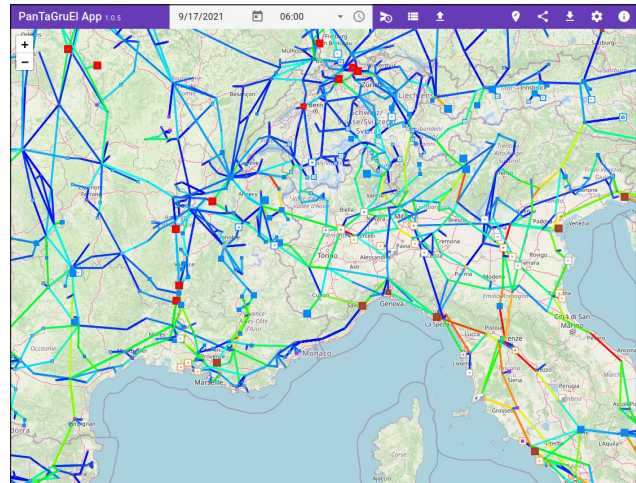
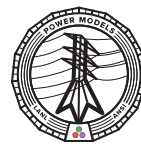
PanTaGruEI

- A Pan-European Transmission Grid and Electricity generation model

- Interactive version at <https://etranselec.ch/pantafrend/>
- L. Pagnier, P. Jacquod, “Inertia location and slow network modes determine disturbance propagation in large-scale power grids”
- M. Tyloo, L. Pagnier, P. Jacquod, “The Key Player Problem in Complex Oscillator Networks and Electric Power Grids: Resistance Centralities Identify Local Vulnerabilities”

- Our version: PowerModels format

- 7822 power lines and 553 transformers
- 4097 buses with load distributed based on population density
- 815 generators of various types



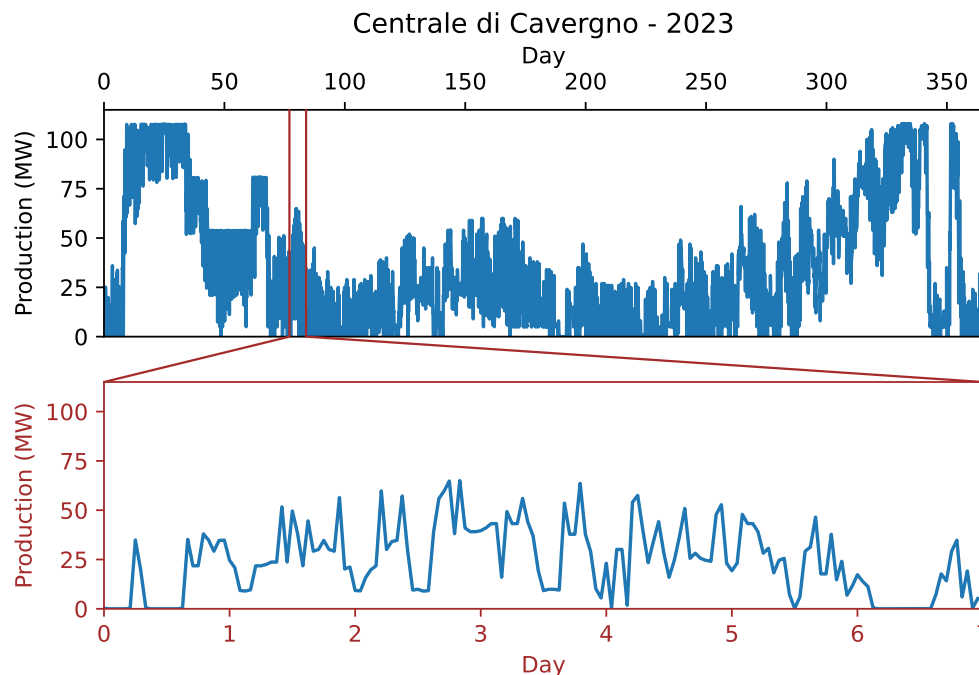
The other problem: time series

- Steady-state data, one hour time resolution: some data available, but not enough for ML...

- Generating synthetic data is very challenging!

- How do we model *this*?

(data source: <https://transparency.entsoe.eu/>)



Typical synthetic data approaches

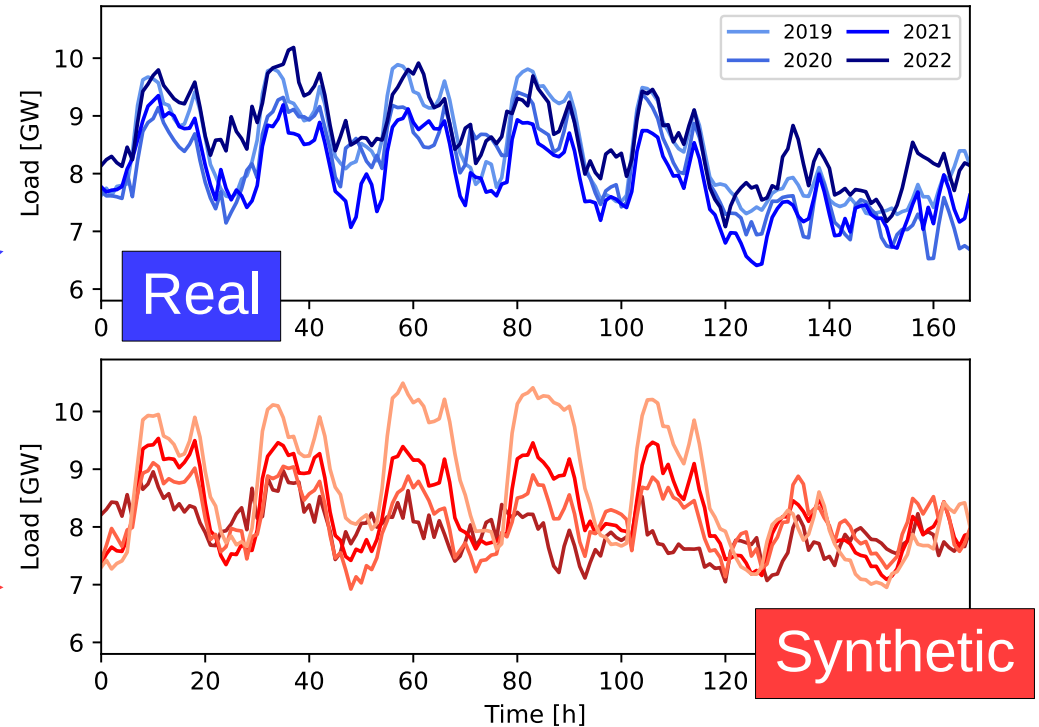
- Often no coherence in time
- When realistic time series exist, typically too few
→ spurious correlations
- Production dispatch from Optimal Power Flow (OPF) leads to unrealistically many saturated lines

Our solution

- Split the problem into load modeling and production dispatch
- **Loads:** build and use a statistical model
- **Production:** optimization problem
 - Thermal limits of the lines as objective instead of hard constraints
 - Constraints on integrated production + ramp constraints
 - Add noise to mimic wild electricity market

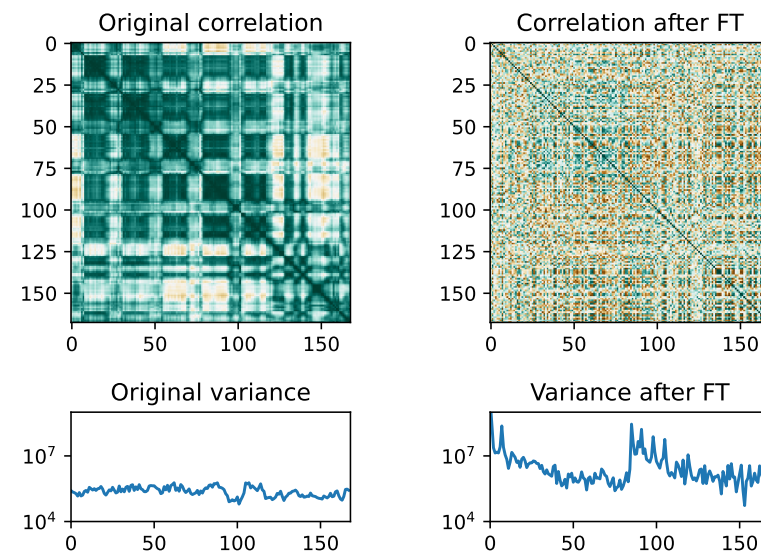
Synthetic load series

- Based on total load by country from public ENTSO-E data <https://transparency.entsoe.eu/>
- Use multiple years to estimate variance
- Build multivariate Gaussian distribution



Technicalities

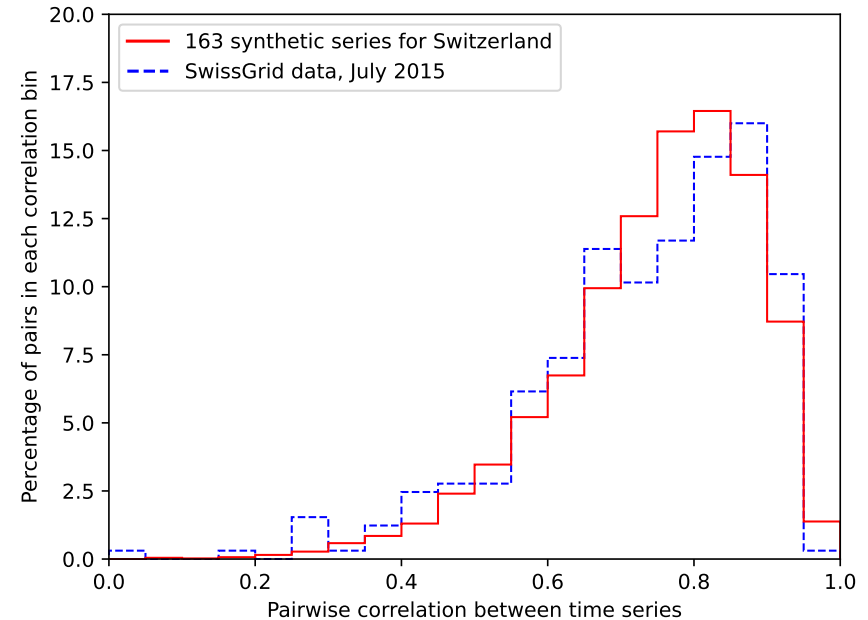
- Long time series → big correlation matrix
- Daily, weekly & yearly patterns → large off-diagonal entries
- Pass through Fourier transform for efficient modeling
 - Few “signal” frequencies, highly correlated
 - Many “noise” components, uncorrelated
 - FFT algorithm



Example with one-week long series

Realistic load series


- Arbitrarily many series can be generated quickly
- Mean corresponding to historical value
- Adjustable variance
→ fit correlations
- Room for improvement with more granular data, other scenarios, ...



Production dispatch

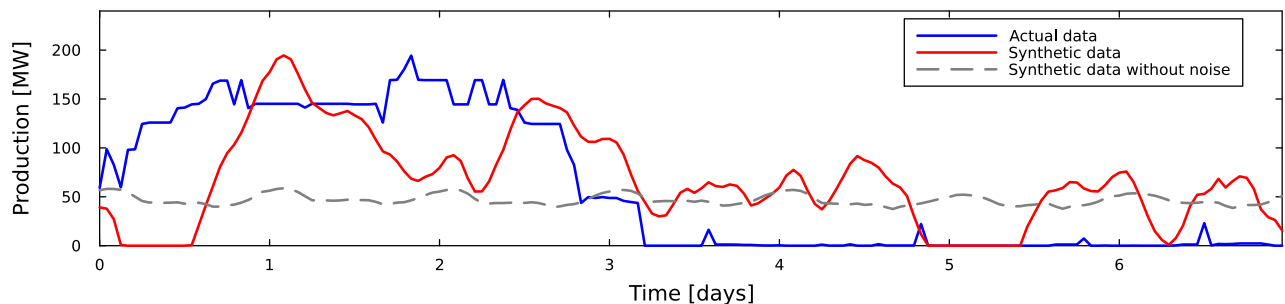
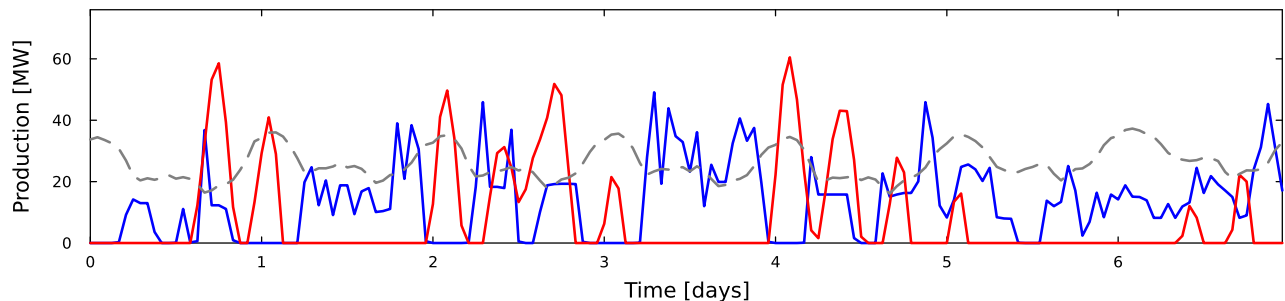
- Optimization problem (~~economic OPF~~)
- Minimize power flows through lines
 - penalize overloaded lines with quadratic cost (instead of hard constraint)
 - achieve relatively local production like TSO
- Annual constraints based on published availability
- Treat non-dispatchable sources separately (nuclear)

Optimal Power Flow

$$\min_{\vec{P}_t} \left[\sum_t \vec{P}_t^T \mathbf{A} \vec{P}_t + \sum_t \vec{b}_t \vec{P}_t \right] \quad \text{s.t.} \quad 0 \leq \vec{P}_t \leq \vec{P}^{\max}, \quad \sum_t \vec{P}_t = \vec{E}, \quad \vec{P}_t \cdot \vec{1} = L_t$$


- Complex behavior resulting from
 - variable loads
 - linear generation cost (noise with typical freq. + harmonics)
- Optimization problem is convex & feasible
 - runs on laptop with Gurobi: 1 year of data in 2-3 hours

OPF results



- Noise chosen as small as possible, but sufficient to trigger on/off behavior
- Several operating modes captured
- Realistic aggregated production by country and type

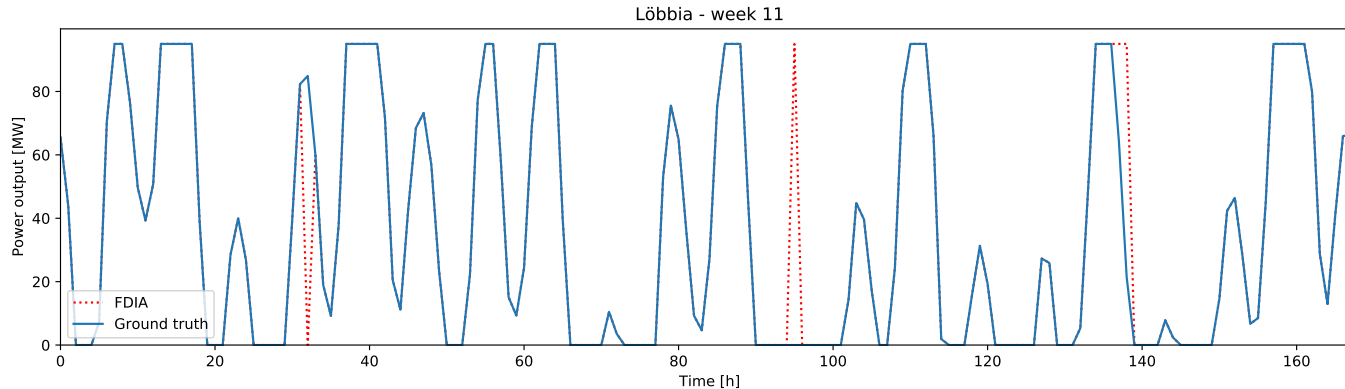
The published dataset

- Data available on Zenodo with data descriptor paper at Nature Scientific Data
- 20 years of time series with one-hour resolution for 815 generators and 4097 loads
- Open-access tools on GitHub repository <https://github.com/GeeeHesso/PowerData>



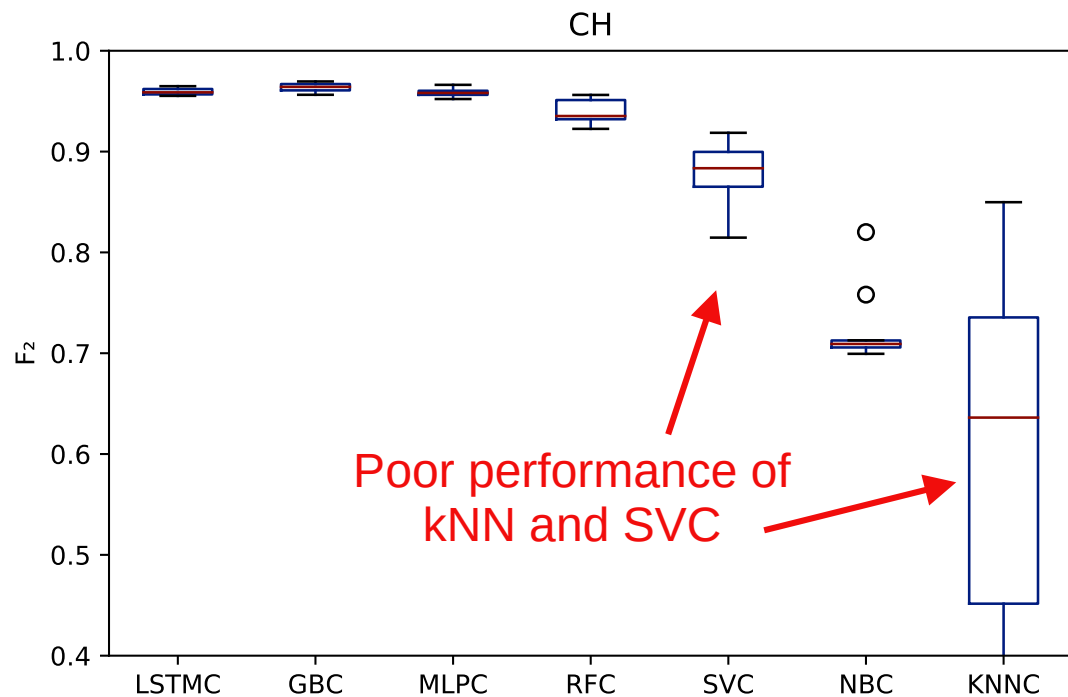
Application: anomaly detection

- Scenario: false data injection attack at one production site

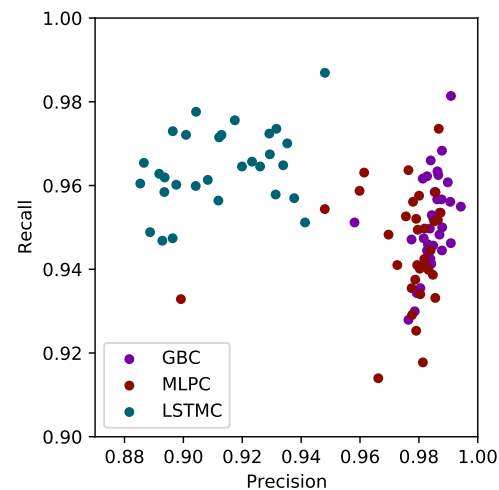


- Impossible to detect in isolation
- Focus on grid of one country (point of view of TSO)
- Use F_2 score (penalize false negatives more than false positives)

Comparing classification algorithms

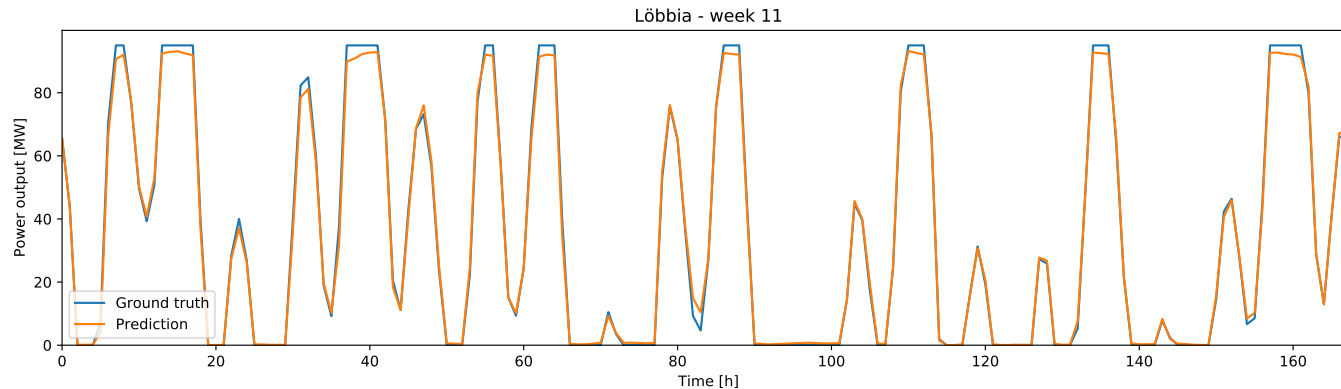


- Good performance of NN (deep and shallow)
- Consistent with purely **contextual** anomalies

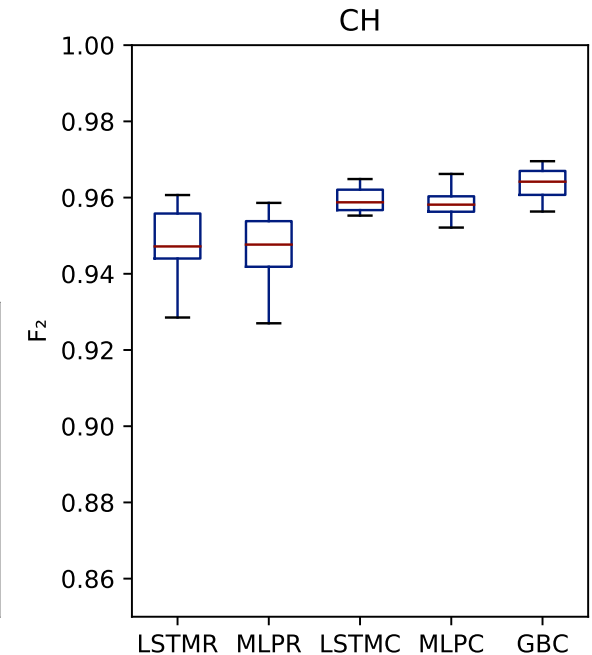


Unsupervised algorithms

- Similar performances with unsupervised algorithms
- Much more versatile

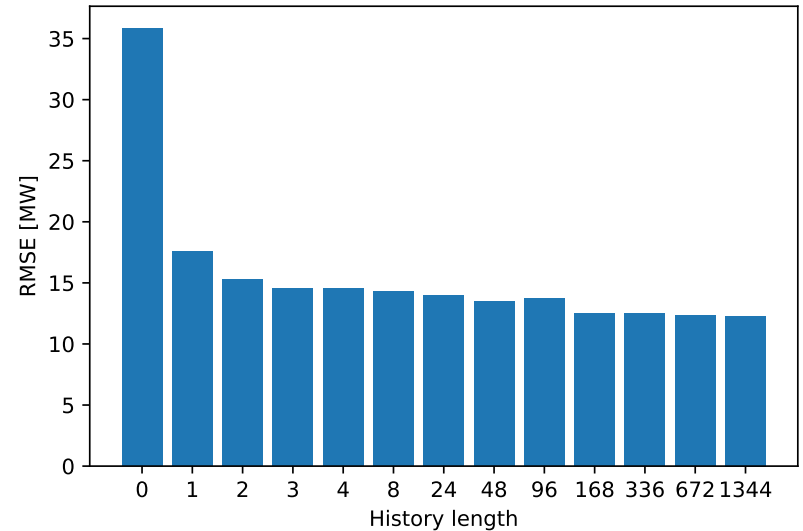


Classification performance



Some lessons for anomaly detection

- History is crucial, but no need to go beyond few hours



- Relevant context: other production sources
- Grid distances have little impact

Conclusions

- A large and realistic **dataset** available to you
- New **methods** for data multiplication & modeling
- **Anomaly detection**: can you beat us with FM?

Thank you!