# DB for AI:
# Data Management for Deep Learning

Lei Chen

**Department of Computer Science and Engineering**

**The Hong Kong University of Science and Technology**

# Outline

- **<u>Background and Motivation</u>**

- Technical Challenges

- Our Reasearch Studies
  - Knowledge Extraction and Labelling
  - Graph substitutions on DNN computation graphs
  - Explainable Recommendation and Explainable GNN

- Beyond DB for AI

# Background: DL applications are ubiquitous

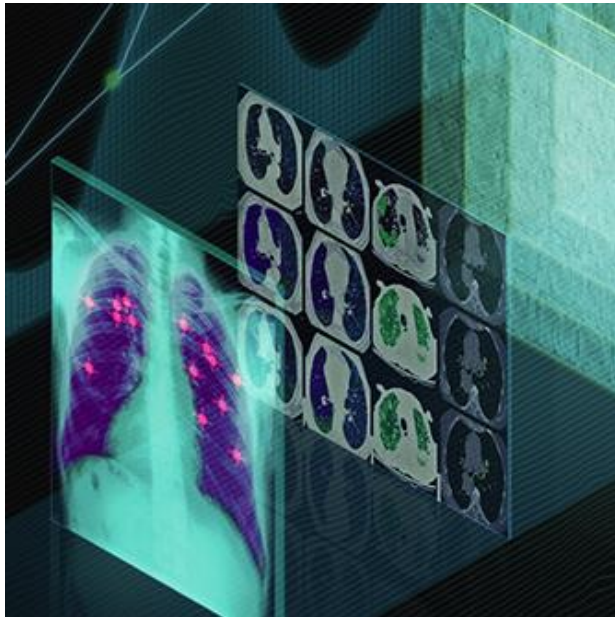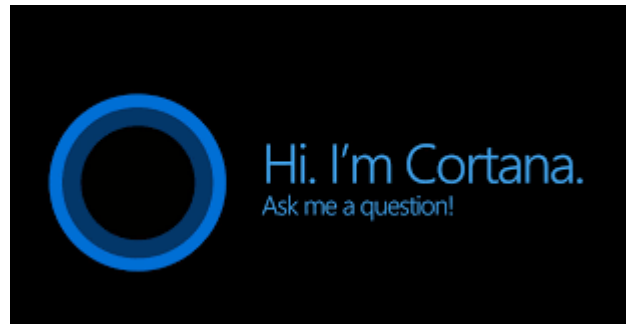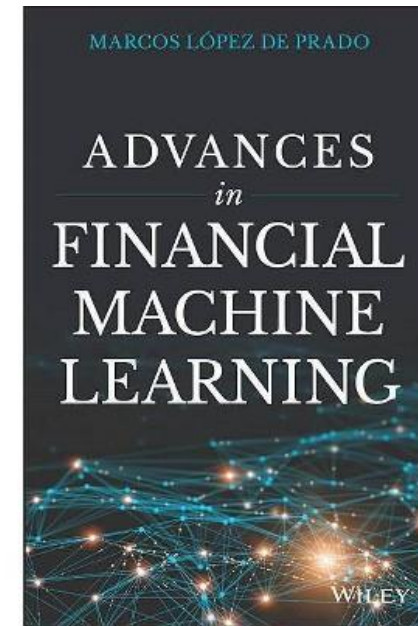- DL has made a huge *success* over the past years.



**Image Recognition**

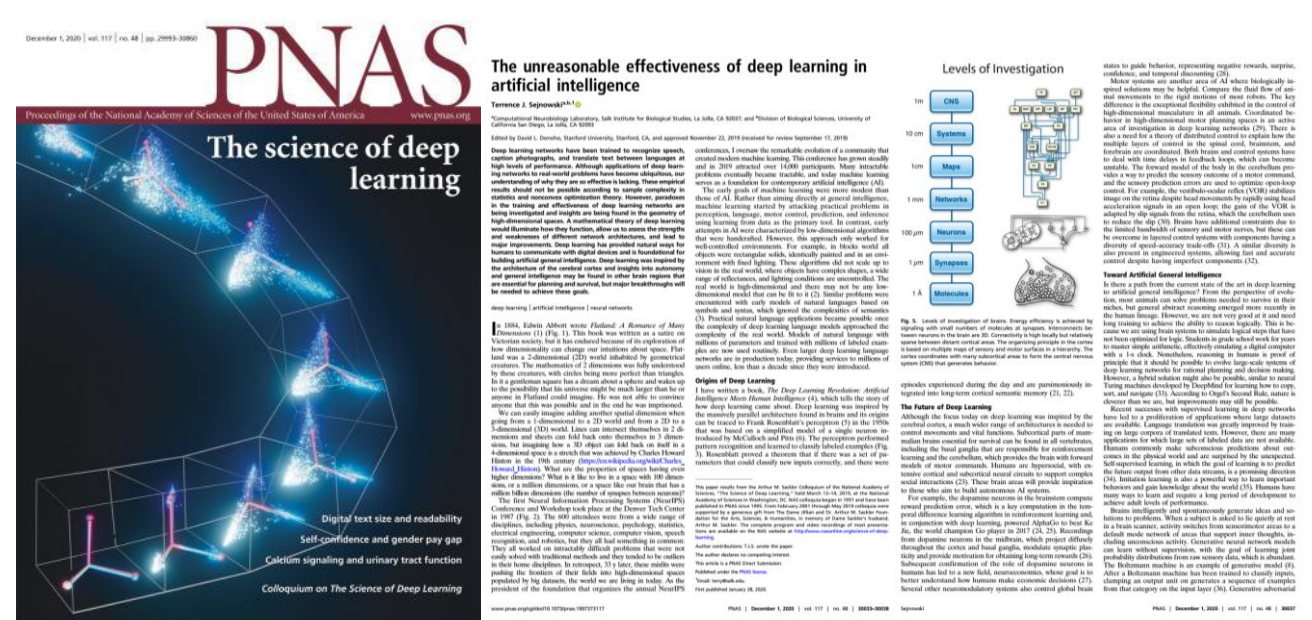**Natural Language Processing**

**Smart Finance**

**Intelligent Transportation**

# Background: Data is the new oil

- The first secret of DL's success: *big data*



"The world's *most valuable resource* is no longer oil, but data".   -- The Economist, 2017

"*Recent successes* in deep networks have led to a proliferation of applications where *large datasets are available*".   -- Terrence J. Sejnowski, in PNAS 2020

# Motivation: Why data management for DL?
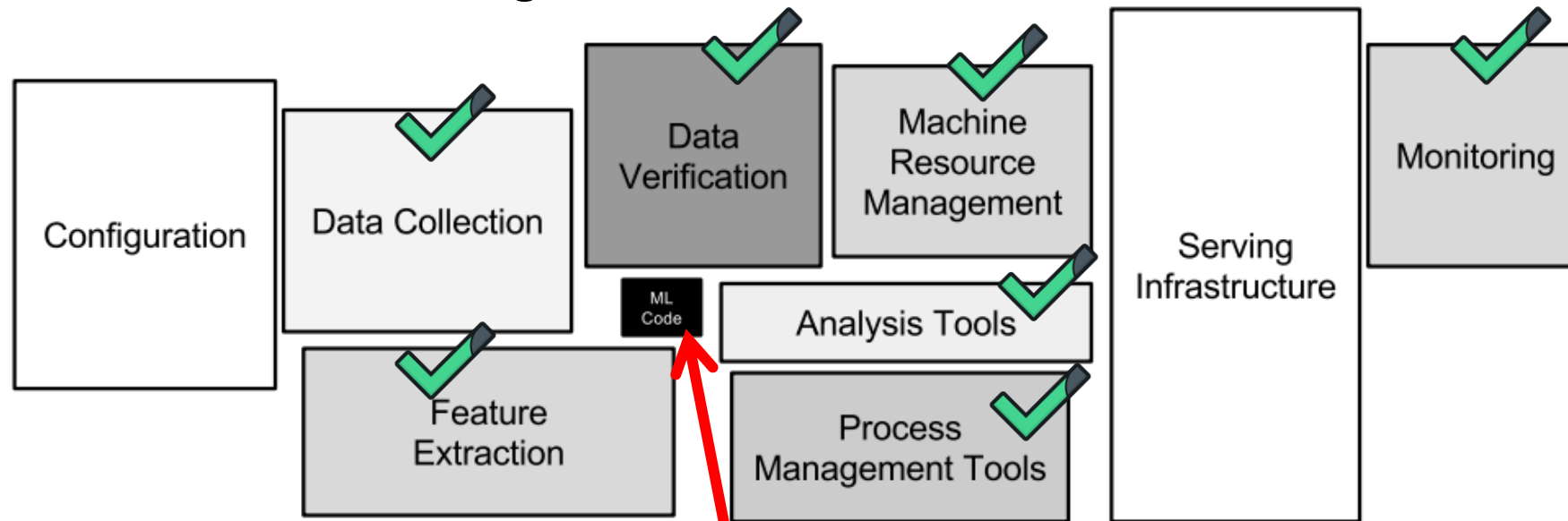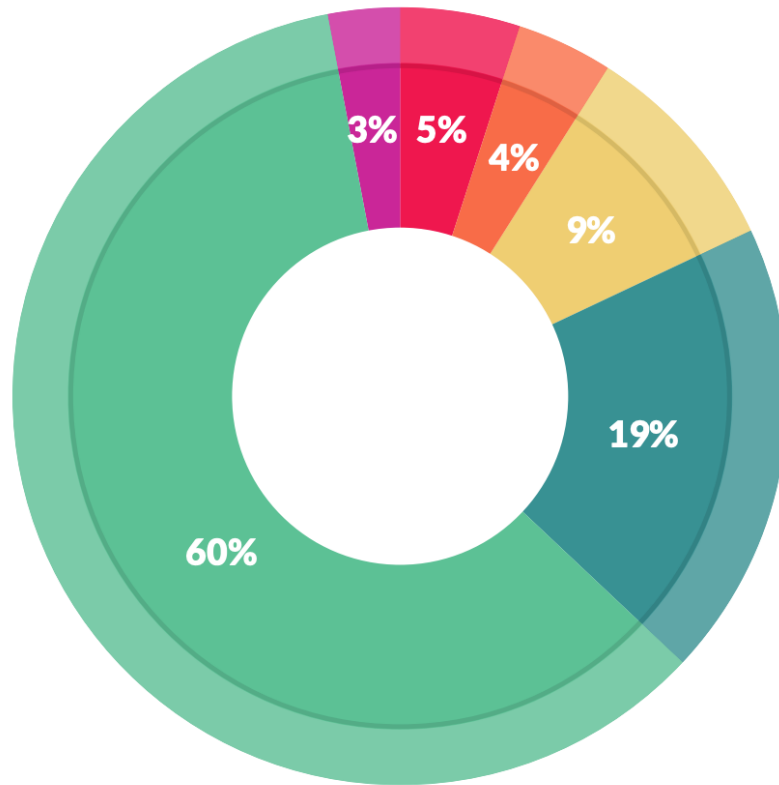
✓ : related to data management



Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

*"In Google, only a tiny fraction of the code in many ML systems is actually devoted to learning."*

D. Sculley, Gary Holt, Daniel Golovin et al. Hidden Technical Debt in Machine Learning Systems. In NeurIPS 2015.
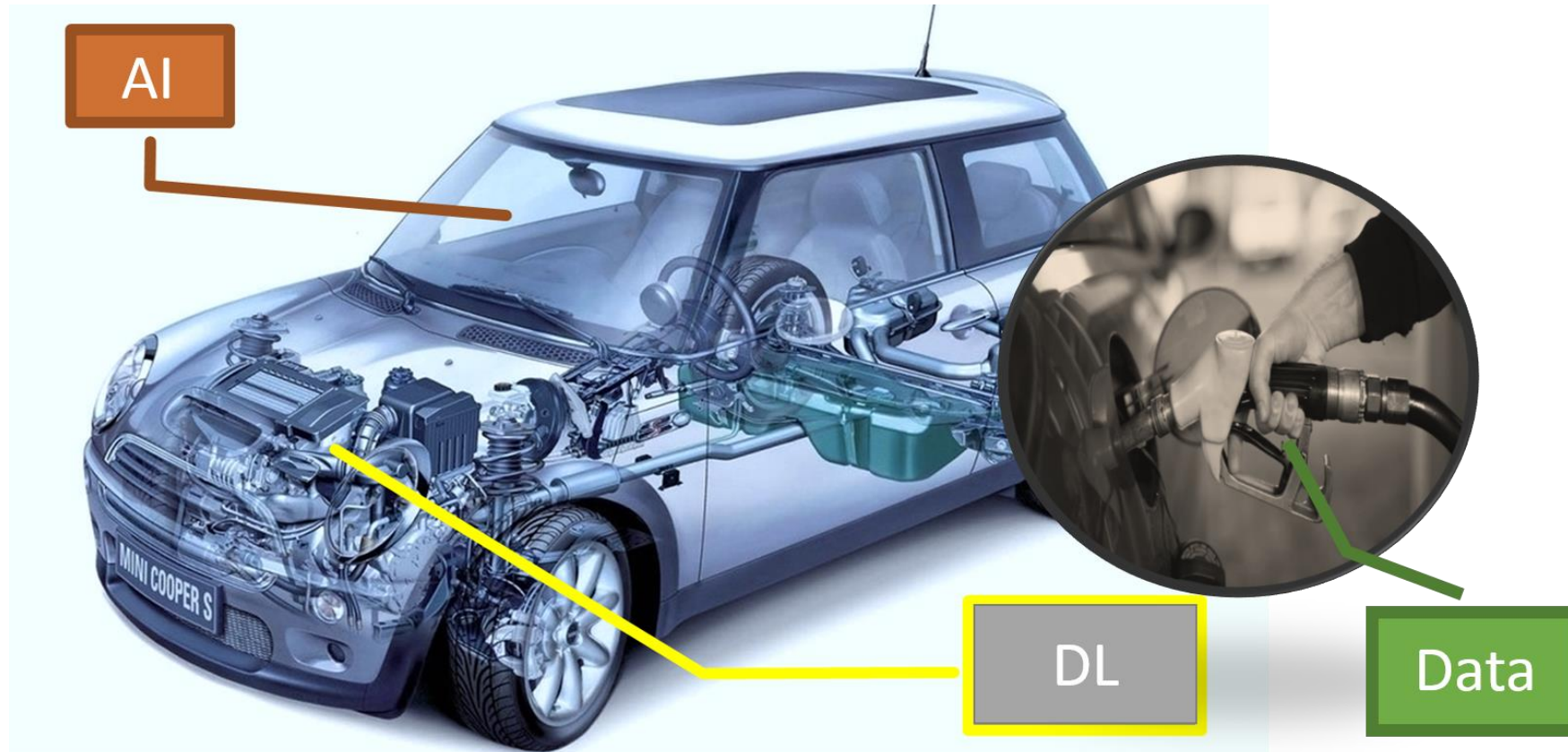
# Motivation: Why data management for DL?



What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

*"80% of ML users' time/effort (often more) spent on data issues!"*

# Background: Big Data, Deep learning and AI

- If data is viewd as the oil, DL is the engine and AI is the car.

# Challenges: data management for DL