# Automatic Machine Learning (AutoML):
# A Tutorial

## Frank Hutter

University of Freiburg
fh@cs.uni-freiburg.de

## Joaquin Vanschoren

Eindhoven University of Technology
j.vanschoren@tue.nl

Slides available at automl.org/events -> AutoML Tutorial
(all references are clickable links)

Computer vision in self-driving cars

Speech recognition







Reasoning in games

**Performance is very sensitive to many hyperparameters**

- Architectural hyperparameters



Units per layer

Kernel size

\# convolutional layers

\# fully connected layers
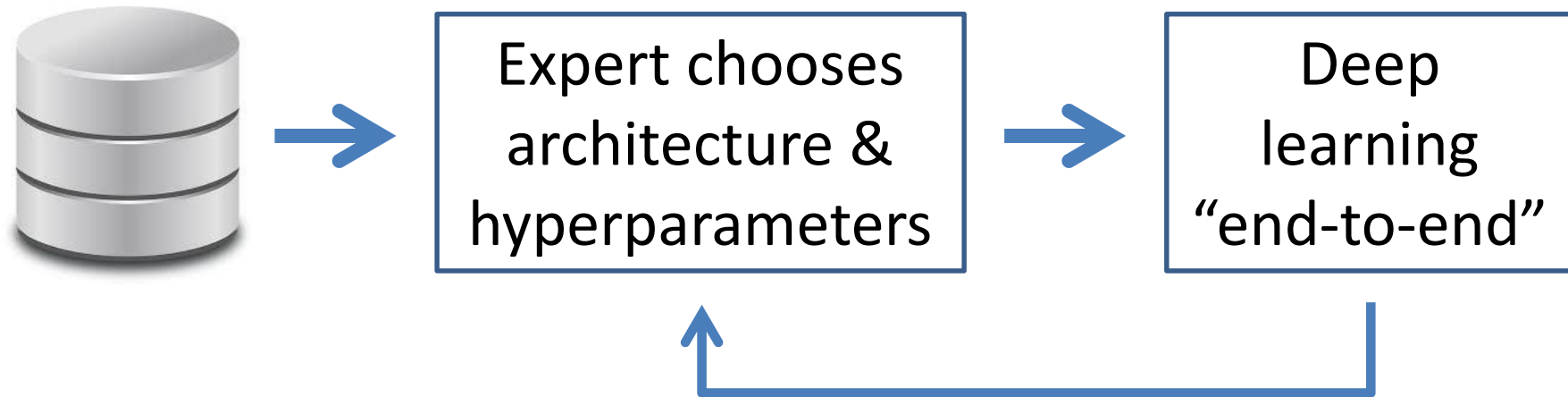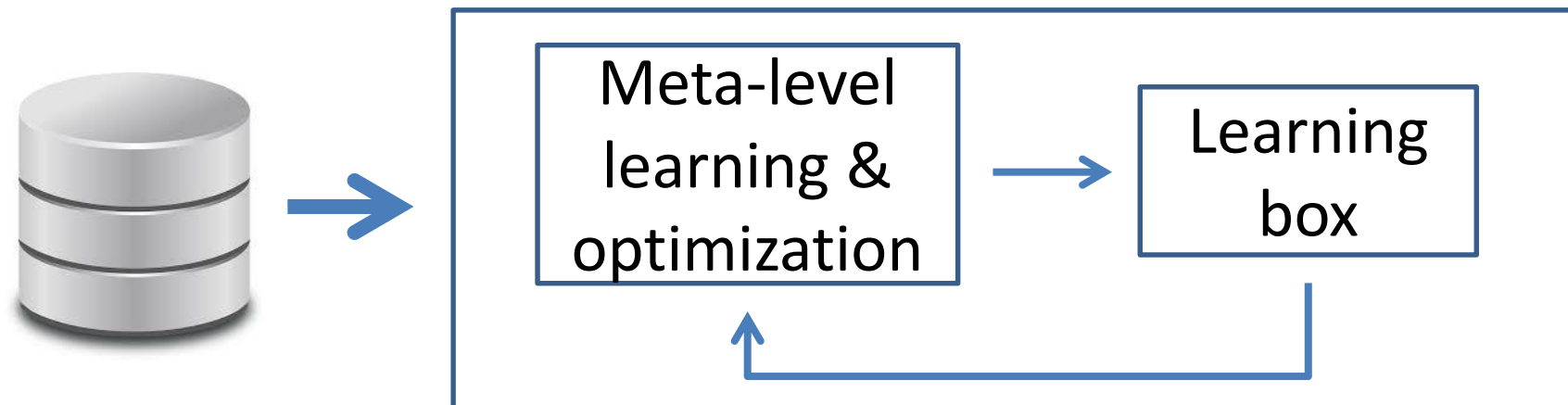
- Optimization algorithm, learning rates, momentum, batch normalization, batch sizes, dropout rates, weight decay, data augmentation, …

→ **Easily 20-50 design decisions**
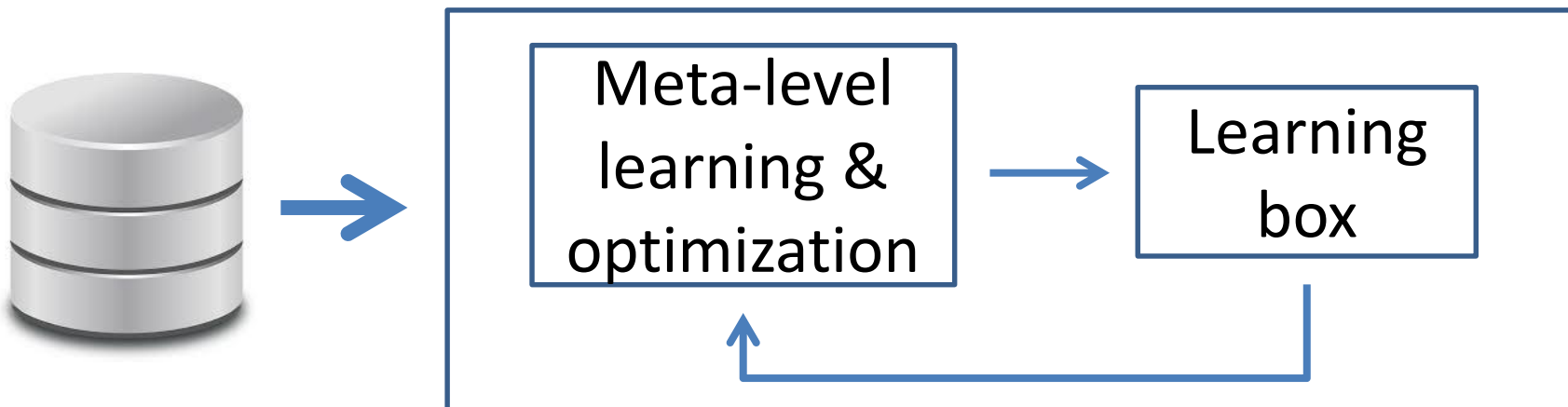
## Current deep learning practice



## AutoML: true end-to-end learning

# Learning box is not restricted to deep learning

- Traditional machine learning pipeline:
  - Clean & preprocess the data
  - Select / engineer better features
  - Select a model family
  - Set the hyperparameters
  - Construct ensembles of models
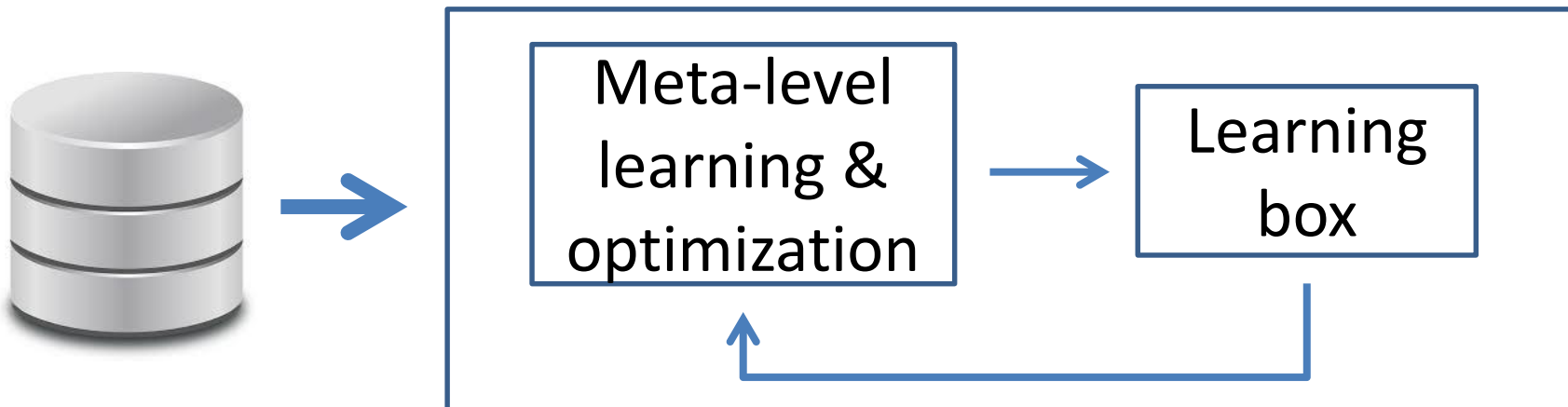  - …

## AutoML: true end-to-end learning

# Outline

1. Modern Hyperparameter Optimization

2. Neural Architecture Search

3. Meta Learning

For more details, see: automl.org/book

## AutoML: true end-to-end learning

# Outline

1. **Modern Hyperparameter Optimization**

   - **AutoML as Hyperparameter Optimization**
   - Blackbox Optimization
   - Beyond Blackbox Optimization

   Based on: Feurer & Hutter: Chapter 1 of the AutoML book: Hyperparameter Optimization

2. **Neural Architecture Search**

   - Search Space Design
   - Blackbox Optimization
   - Beyond Blackbox Optimization

# Hyperparameter Optimization

## Definition: Hyperparameter Optimization (HPO)

Let

- $\boldsymbol{\lambda}$ be the hyperparameters of a ML algorithm $A$ with domain $\boldsymbol{\Lambda}$,
- $\mathcal{L}(A_{\boldsymbol{\lambda}}, D_{train}, D_{valid})$ denote the loss of $A$, using hyperparameters $\boldsymbol{\lambda}$ trained on $D_{train}$ and evaluated on $D_{valid}$.

The hyperparameter optimization (HPO) problem is to find a hyperparameter configuration $\boldsymbol{\lambda}^*$ that minimizes this loss:

$$\boldsymbol{\lambda}^* \in \arg\min_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \mathcal{L}(A_{\boldsymbol{\lambda}}, D_{train}, D_{valid})$$

# Types of Hyperparameters

- **Continuous**
  - Example: learning rate

- **Integer**
  - Example: #units

- **Categorical**

  - Finite domain, unordered
    - Example 1: algo $\in$ {SVM, RF, NN}
    - Example 2: activation function $\in$ {ReLU, Leaky ReLU, tanh}
    - Example 3: operator $\in$ {conv3x3, separable conv3x3, max pool, …}
  - Special case: binary

# Conditional hyperparameters

- **Conditional hyperparameters** B are only active if other hyperparameters A are set a certain way

  - Example 1:
    - A = choice of optimizer (Adam or SGD)
    - B = Adam's second momentum hyperparameter (only active if A=Adam)

  - Example 2:
    - A = type of layer k (convolution, max pooling, fully connected, …)
    - B = conv. kernel size of that layer (only active if A = convolution)

  - Example 3:
    - A = choice of classifier (RF or SVM)
    - B = SVM's kernel parameter (only active if A = SVM)

# AutoML as Hyperparameter Optimization

## Definition: Combined Algorithm Selection and Hyperparameter Optimization (CASH)

Let

- $\mathcal{A} = \{A^{(1)}, \ldots, A^{(n)}\}$ be a set of algorithms
- $\mathbf{\Lambda}^{(i)}$ denote the hyperparameter space of $A^{(i)}$, for $i = 1, \ldots, n$
- $\mathcal{L}(A_{\boldsymbol{\lambda}}^{(i)}, D_{train}, D_{valid})$ denote the loss of $A^{(i)}$, using $\lambda \in \mathbf{\Lambda}^{(i)}$ trained on $D_{train}$ and evaluated on $D_{valid}$.

The Combined Algorithm Selection and Hyperparameter Optimization (CASH) problem is to find a combination of algorithm $A^* = A^{(i)}$ and hyperparameter configuration $\boldsymbol{\lambda}^* \in \mathbf{\Lambda}^{(i)}$ that minimizes this loss:

$$A_{\boldsymbol{\lambda}^*}^* \in \underset{A^{(i)} \in \mathcal{A}, \boldsymbol{\lambda} \in \mathbf{\Lambda}^{(i)}}{\arg\min} \mathcal{L}(A_{\boldsymbol{\lambda}}^{(i)}, D_{train}, D_{valid})$$

→ Simply a HPO problem with a top-level hyperparameter (choice of algorithm) that all other hyperparameters are conditional on
  - E.g., Auto-WEKA: 768 hyperparameters, 4 levels of conditionality