# PART I :

# MACHINE LEARNING

# POPULAR SPECIFICATIONS

* examples :   Inductive program synthesis
    $\equiv$ Programming by Example

  ⊕ Clean notion of a solution program

  ⊖ Very partial specification: many solution programs

* natural language specification :

  ⊕ easy to write

  ⊖ ambiguous: no notion of a solution program

What we'll say here applies to both

# FIRST IDEAS

- 1989, Solomonoff: use specification for creating statistics
- 2013, Menon et al: learn statistics from specification

# DEEP LEARNING

- 2017, Balog et al: DeepCoder neural networks make predictions

# LLMs

- 2021, Open AI: "If we can translate English to French, then we can do code generation: translate language to code!"
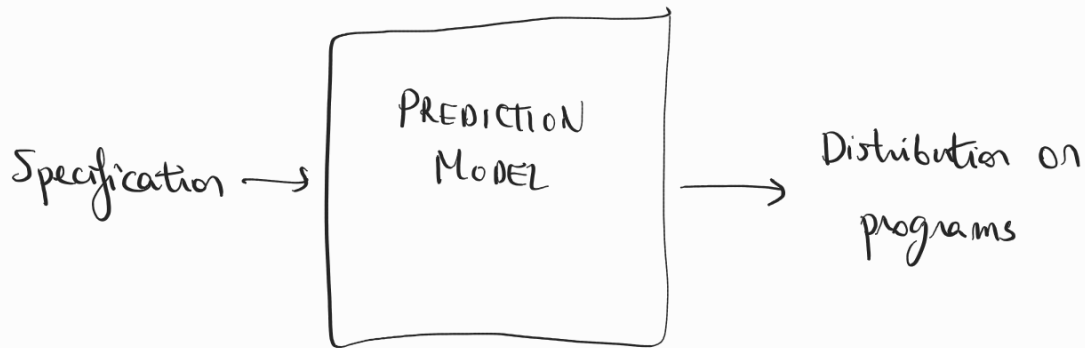
# AN EXAMPLE

Examples:

$$\text{Input} \longrightarrow \text{Output}$$

$$[1, 5, 4, 2) \longrightarrow [2, 4]$$
$$[6, 3, 0, 8] \longrightarrow [0, 6, 8]$$

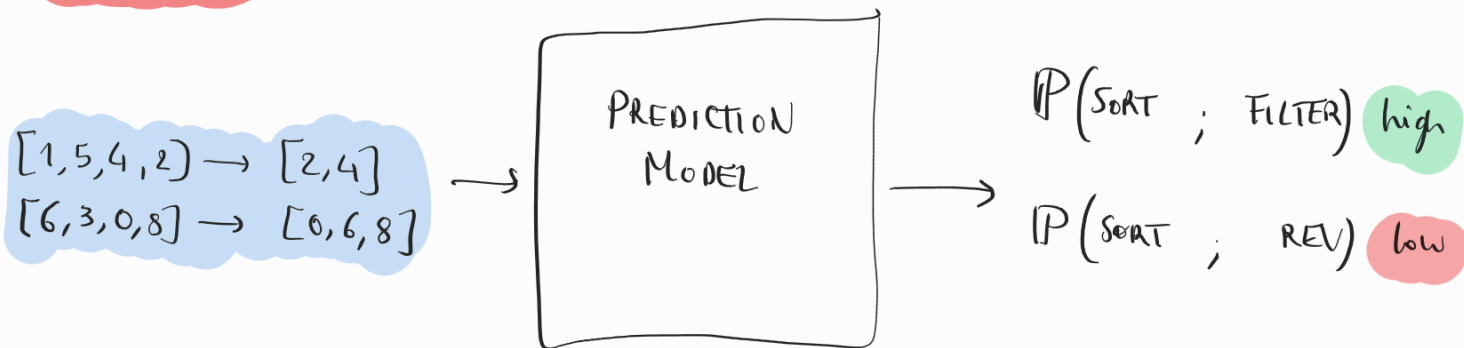Key idea  patterns found in specification reveal which constructs are used in a solution program

Patterns:
- all values in output are even
  - outputs are sorted
  - values in outputs appear in inputs
  ⋮

# PREDICTIONS

Specification $\longrightarrow$ | PREDICTION MODEL | $\longrightarrow$ Distribution on programs

Loss function: $\underset{\text{parameters}}{\max}$ Probability(Program $\models$ Specification)

Example:

$$[1,5,4,2) \longrightarrow [2,4]$$
$$[6,3,0,8] \longrightarrow [0,6,8]$$

$\longrightarrow$ | PREDICTION MODEL | $\longrightarrow$

$\mathbb{P}(\text{SORT} \; ; \; \text{FILTER})$ high

$\mathbb{P}(\text{SORT} \; ; \; \text{REV})$ low

The prediction model intuces:

$$\mathcal{D}(\text{program} \mid \text{specification})$$

the larger $\mathcal{D}(\text{Prog} \mid \text{Spec})$ the more likely $\text{Prog} \models \text{Spec}$
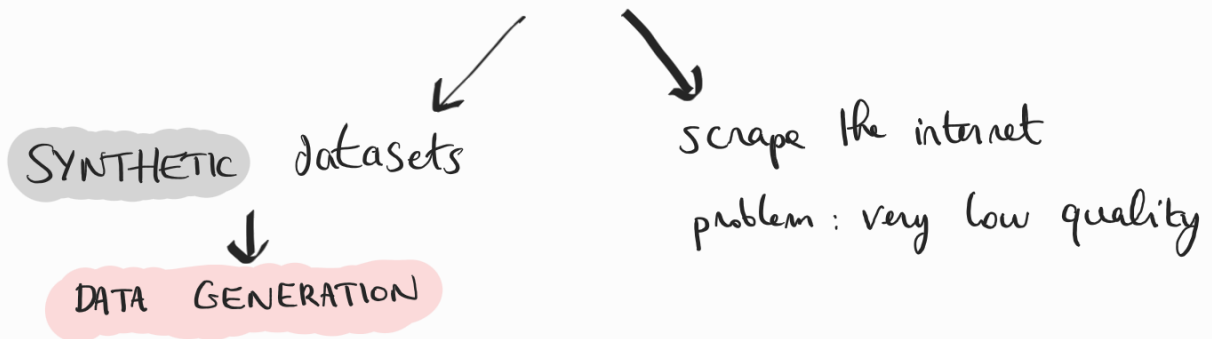
# TRAINING

- **Find** a dataset of $(Spec_i, Prog_i)_{i \in [1,N]}$

- Train the prediction model to maximise

$$\sum_{i=1}^{N} \mathcal{D}(Prog_i \mid Spec_i)$$

# TWO ISSUES

- In practice hard to FIND datasets:

**SYNTHETIC** datasets

$\downarrow$

**DATA GENERATION**

scrape the internet

problem: very low quality

- Program aliasing makes learning harder:

$$\begin{cases} Prog_1 \models Spec \\ Prog_2 \models Spec \end{cases} \Rightarrow \begin{array}{l} \mathcal{D}(prog_1 \mid Spec) = 1/2 \\ \mathcal{D}(prog_2 \mid Spec) = 1/2 \end{array}$$

$\rightarrow$ "programming style" of LLMs

# DATA GENERATION

Step 1 :  Program generation

Step 2 :  For each program, Specification generation

Both are non trivial :

- covering criteria ?  (≠ testing)
- addressing program aliasing ?

A WIDELY OPEN PROBLEM !