# Applied Generative AI

Weak Supervision for Improved Text-to-Label Tasks
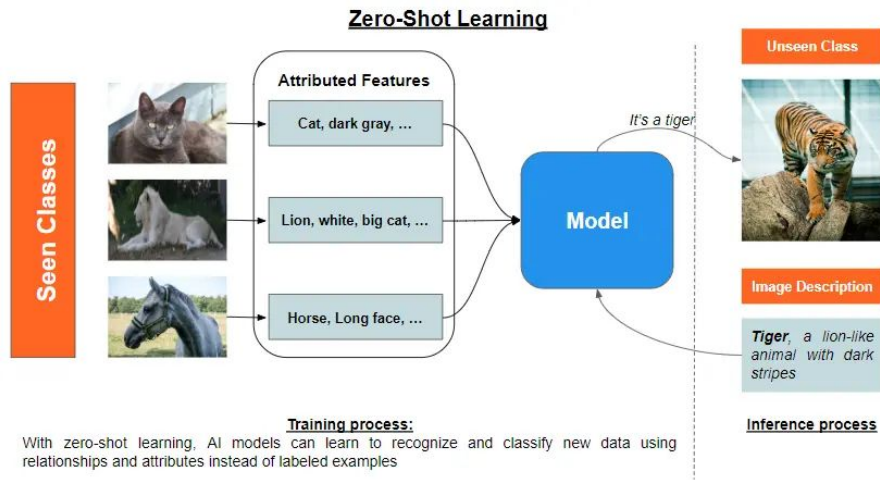
# The Labeled Data Bottleneck



- Quality, labeled data remains a frequent bottleneck for ML application.

- Data quality frequently determines ML project success.

- Getting quality, labeled data is frequently resource and time intensive.

# The Labeled Data Bottleneck



**Zero-Shot Learning**

**Attributed Features**

Cat, dark gray, ...

Lion, white, big cat, ...

Horse, Long face, ...

**Model**

Seen Classes

**Unseen Class**

It's a tiger

**Image Description**

*Tiger*, a lion-like animal with dark stripes

**Inference process**

**Training process:**
With zero-shot learning, AI models can learn to recognize and classify new data using relationships and attributes instead of labeled examples
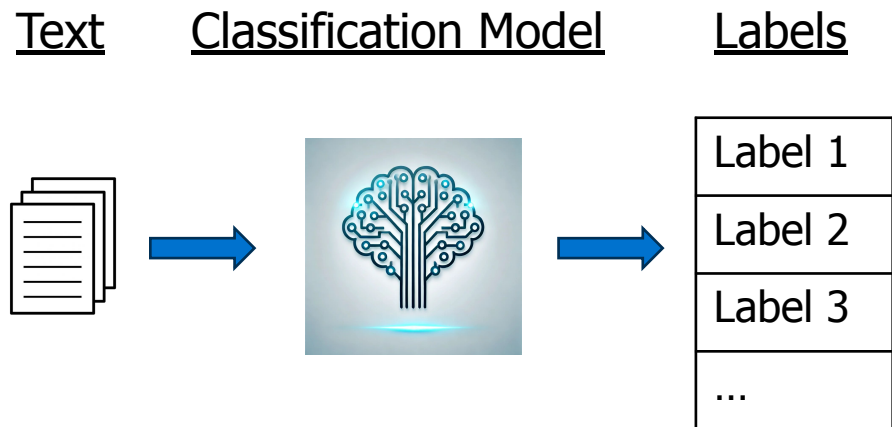
Vina, Abirami. *Understanding few-shot, zero-shot, and transfer learning*. Ultralytics Blog. 2025

- Several research fields attempt to address this problem
  - Zero-shot learning
  - Prompt Engineering
  - And many more...
- Recent emergent properties in large models are the primary method
  - Still frequently worse that in-domain models

# Agenda

- Review of Text-to-Label and Prompt Engineering
- Weak Supervision
- Combining GAI and Weak Supervision
- Code Example

# **Text-to-Label Task Review**

Text         Classification Model         Labels



| Label 1 |
| Label 2 |
| Label 3 |
| ... |

- Tasks involving classification or label assignment to text inputs.
- Examples:
  - Sentiment analysis
  - Topic classification
  - Spam detection
  - Code bugs

# Text-to-Label with Generative AI
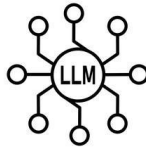


Prediction/ Inference

Text to be labeled

LLM

Labeling Prompt

Examples

Post- Processing

Labels    Explanations

- Labeling by Generative AI
1. Develop prompting scheme
2. Add text to be labeled to prompt and send model
3. Post-process output
4. Refine the prompt based on some labeled examples

- Consider few-shot examples and batch prompting, if applicable

# Prompt Engineering for Text-to-Label

```
[prompt] Stance classification
is the task of understanding a
person's opinion, either implied
or expressed, toward a target.
Classify the stance of the
statement below toward the
target below.
Target: {target}
Text: {text}
Stance:
```
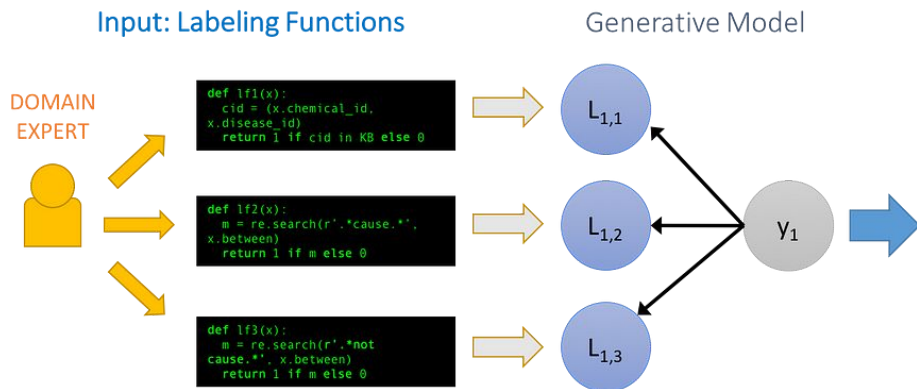
- *Prompt Engineering* is the art and science of designing and structuring *prompts* (questions or tasks) fed to language models.

- When doing text-to-label, clear task instructions, definitions, and indicators are usually very important.

- Few-shot Prompting and Chain-of-Thought are frequently used patterns for Text-to-Label.

# What is Weak Supervision?



Input: Labeling Functions

Generative Model

DOMAIN
EXPERT

```
def lf1(x):
    cid = (x.chemical_id,
    x.disease_id)
    return 1 if cid in KB else 0
```

$L_{1,1}$

```
def lf2(x):
    m = re.search(r'.*cause.*',
    x.between)
    return 1 if m else 0
```

$L_{1,2}$

$y_1$

```
def lf3(x):
    m = re.search(r'.*not
    cause.*', x.between)
    return 1 if m else 0
```

$L_{1,3}$
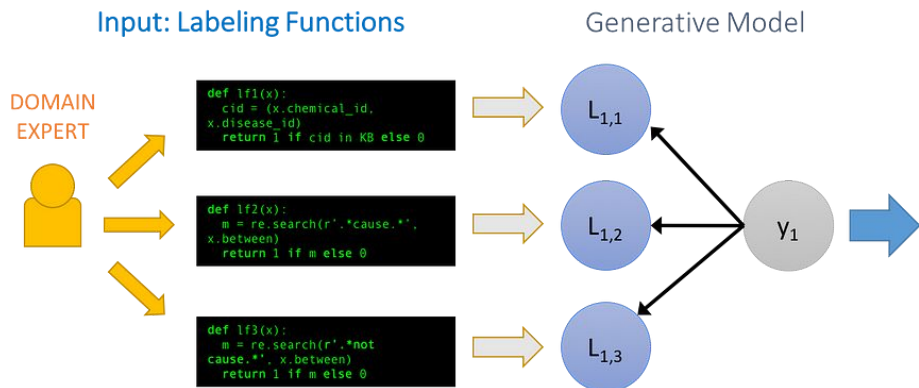
SAIL Blog, *Weak Supervision: A New Programming Paradigm for Machine Learning* (2019)

- A technique for labeling data using noisy, incomplete, or imprecise sources.
  - Combines multiple weak signals to create higher-quality labels.

- Popular package for weak supervision is Snorkel
  - Consists of labeling functions and a generative labeling model.
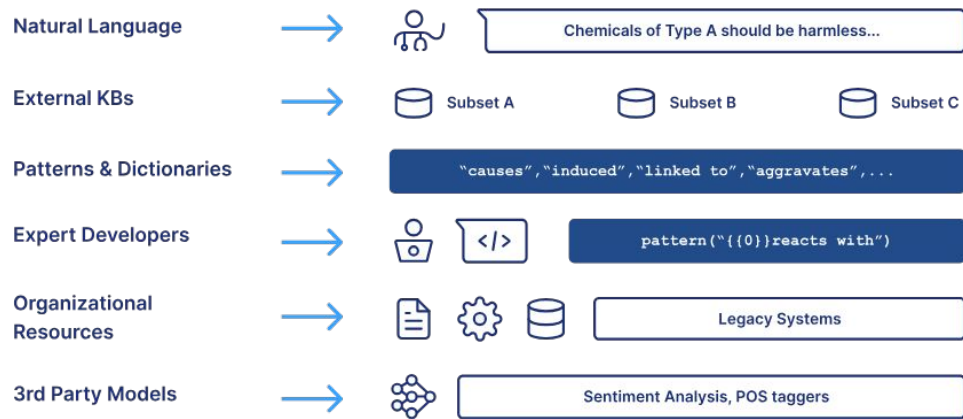
# Using Weak Supervision in Snorkel



SAIL Blog, *Weak Supervision: A New Programming Paradigm for Machine Learning* (2019)

- Objective is to combine multiple weak signals of the label to create higher-quality labels.

- For snorkel, we need to define the weak labels by "Labeling Functions"

- After the producing the weak labels, we can evaluate them and combine them with a generative model to produce the quality labels
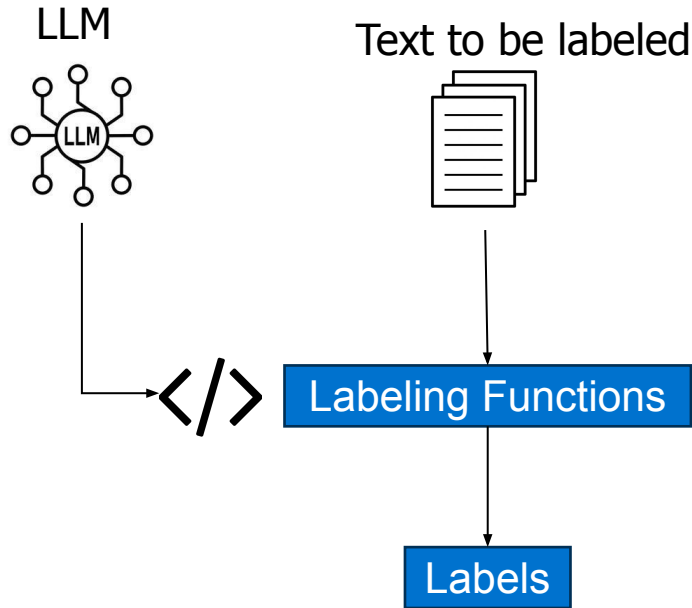
# Creating Labeling Functions



Snorkel Team, *Weak Supervision* (2023)

- Take in an example and produce a weak label.
  - Weak label can also be an "abstain" or no-label

- Functions are normal Python functions marked with labeling function decorator

- Can basically use just about anything (and any kind of signal) to create labeling functions

# Creating Labeling Functions by LLM

LLM

Text to be labeled

Labeling Functions

Labels

- First proposed by Huang et al. in 2024, ALCHEmist approach uses a LLM to create labeling functions.

- Much more scalable approach than direct labeling

- Typically works well only when combined with weak supervision

# Key Takeaways

- Weak supervision is a great way to turbo-charge your GAI derived labels when creating text-to-label solutions

- Weak supervision is a great way to combine various types of signals to create quality data labels

- Typically, you want to combine the final labels from weak supervision with a lot of data and training a model to then iterate toward even higher quality labels and models