

Composable Machine Learning

Eric Xing
Petuum & Carnegie Mellon



Acknowledgements



Zhiting Hu

Scientist @Petuum Inc.
PhD Candidate @CMU



Qirong Ho

Co-founder, CTO @Petuum Inc.

R&D team @

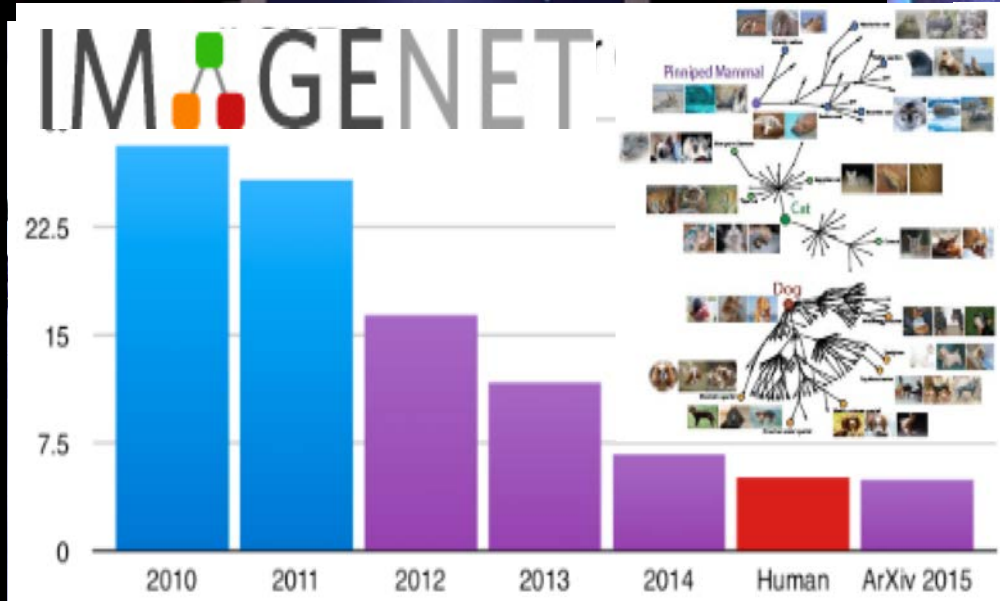
Petuum



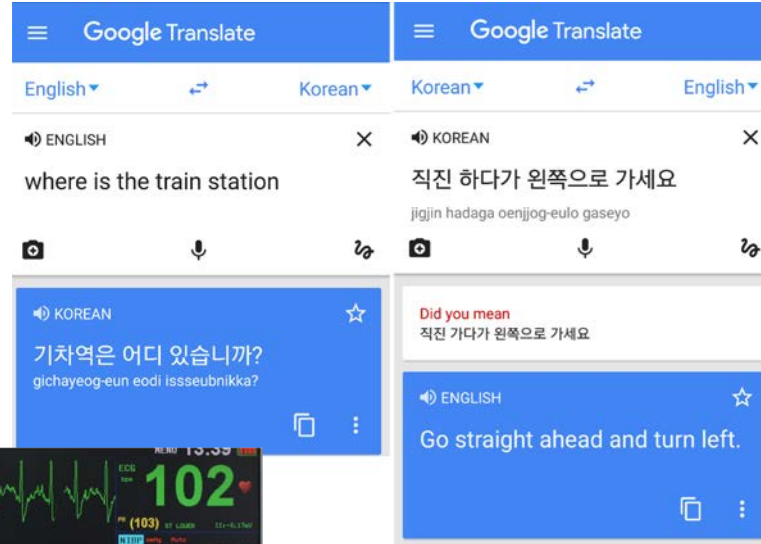
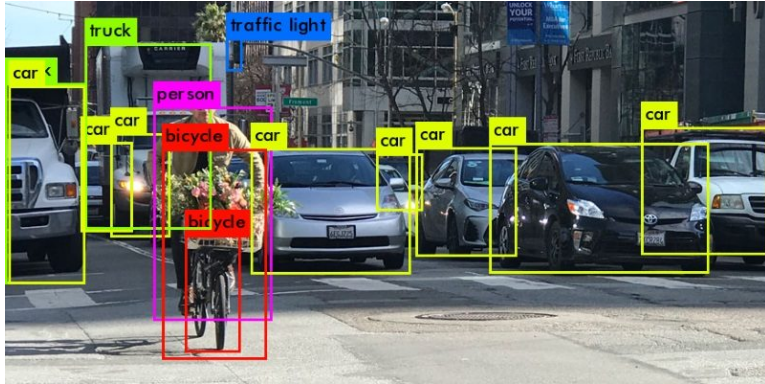
Sailing Lab @

**Carnegie
Mellon
University**

An AI era?



Real-world Machine Learning Problems



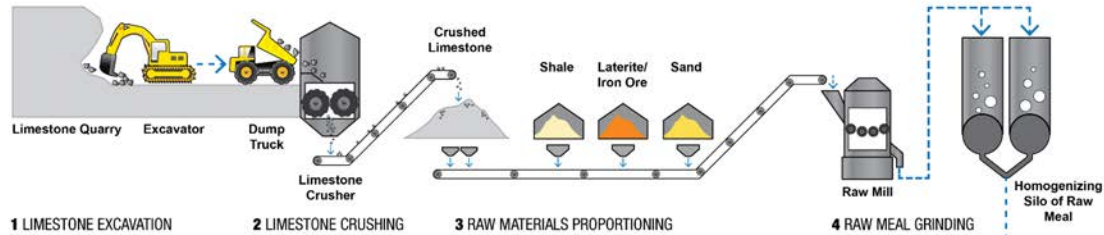
The Cement Industry

PRODUCTION PROCESS OF CEMENT



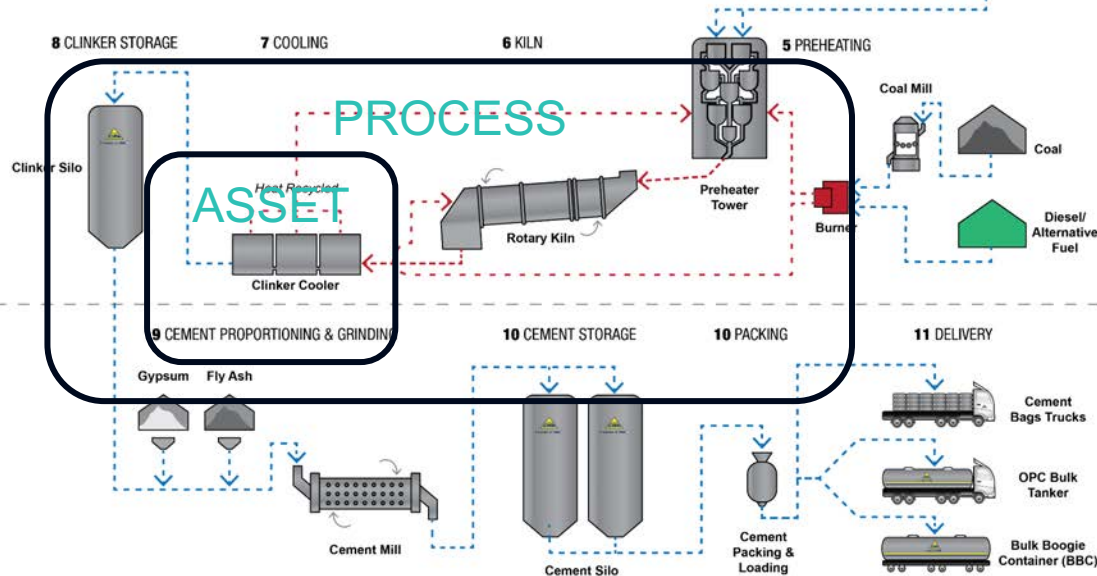
STAGE 1 RAW MATERIALS PREPARATION

Raw materials needed to produce cement i.e. calcium carbonate, silico, alumina and iron ore are extracted from limestone rock, chalk, shale or clay and ferum containing material. These raw materials are crushed through a milling process.



STAGE 2 PYRO PROCESSING

Grinding produces a fine powder, known as raw meal, which is preheated and then sent to the kiln. The raw meal is heated to around 1,450 °C, where chemical reactions take place to form cement clinker.



STAGE 3 CEMENT GRINDING AND DISTRIBUTION

A small amount of gypsum is added to the clinker to regulate cement setting time. The mixture is then very finely ground to obtain 'pure cement'. The cement is stored in silos or storages before being distributed in bulk or in bags to the sites where it will be used.

Operational Excellence Initiatives

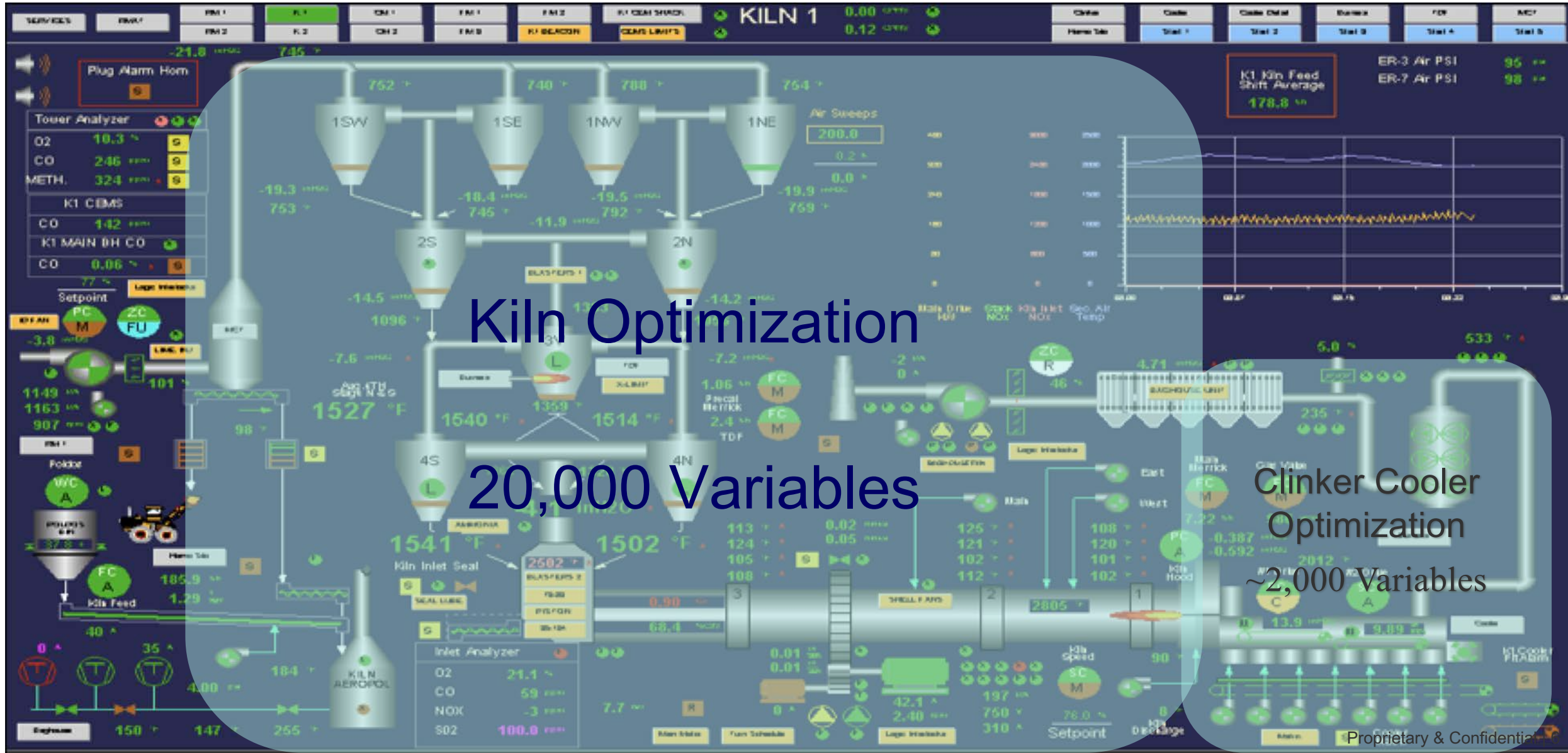
Emissions (NO_x, SO_x)

Fuel-Mix (Traditional vs. Renewables vs. Alternatives)

Benchmarking fleet Performance

Predictive/Preventive Maintenance

Assess, Process, Factory, Corporate Optimization



Kiln Optimization

20,000 Variables

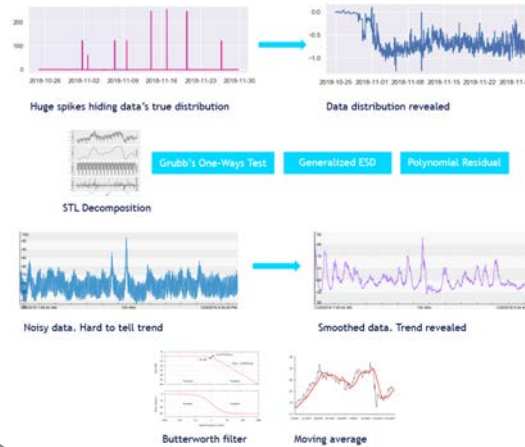
Clinker Cooler Optimization
~2,000 Variables

Building an ready-to-use AI solution for this is

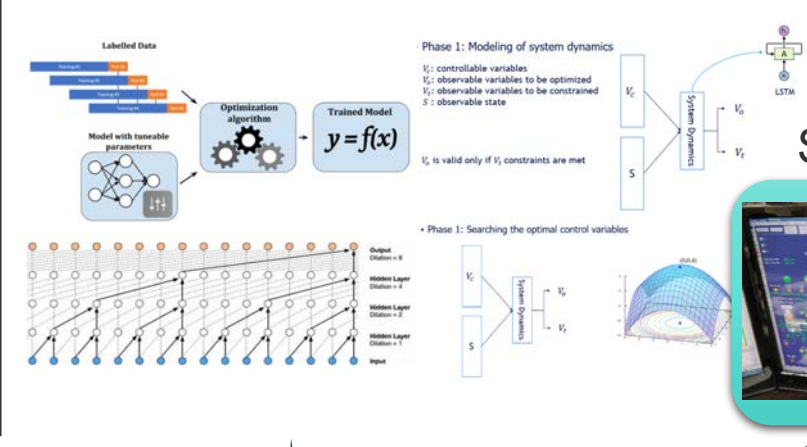
Extremely complex



Raw data far from ideal



Simultaneous Prediction and Control



Systems/Infra



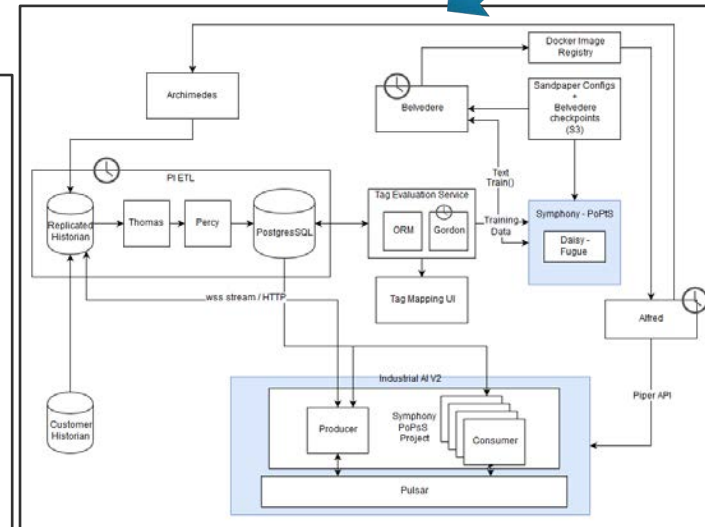
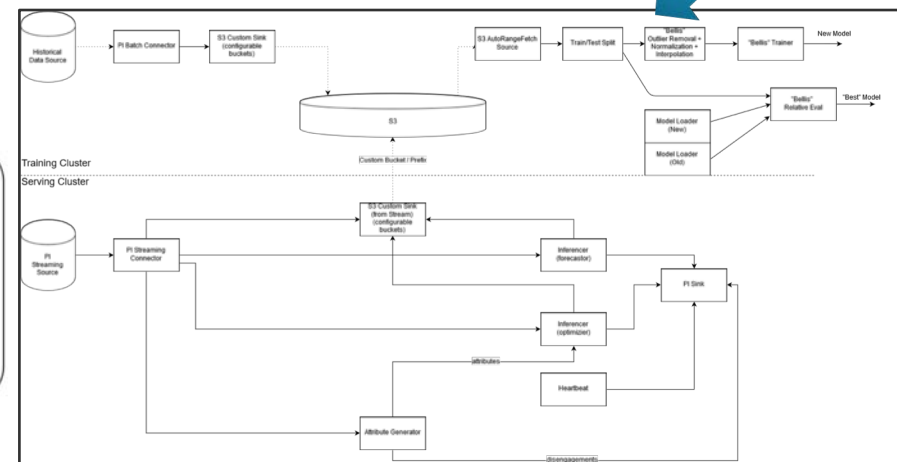
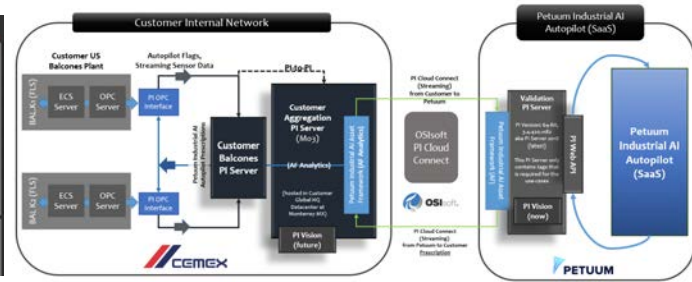
Requires inter-plant operation between diverse systems

Can't solve with "algo marketplace" or Kaggle competition

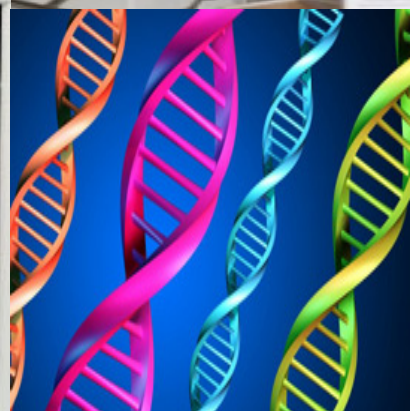
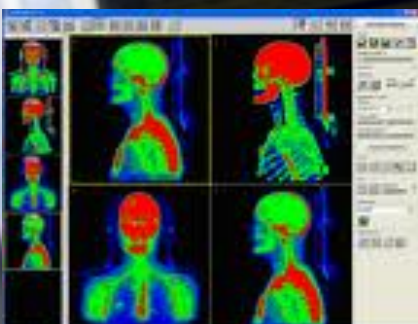
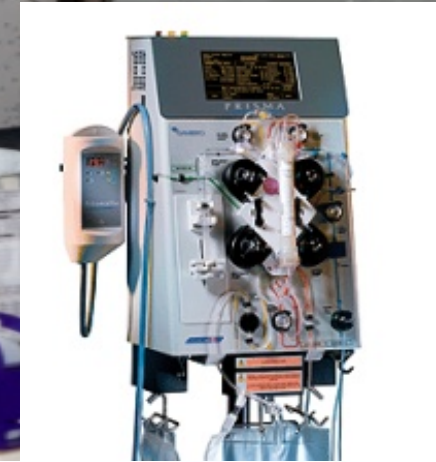
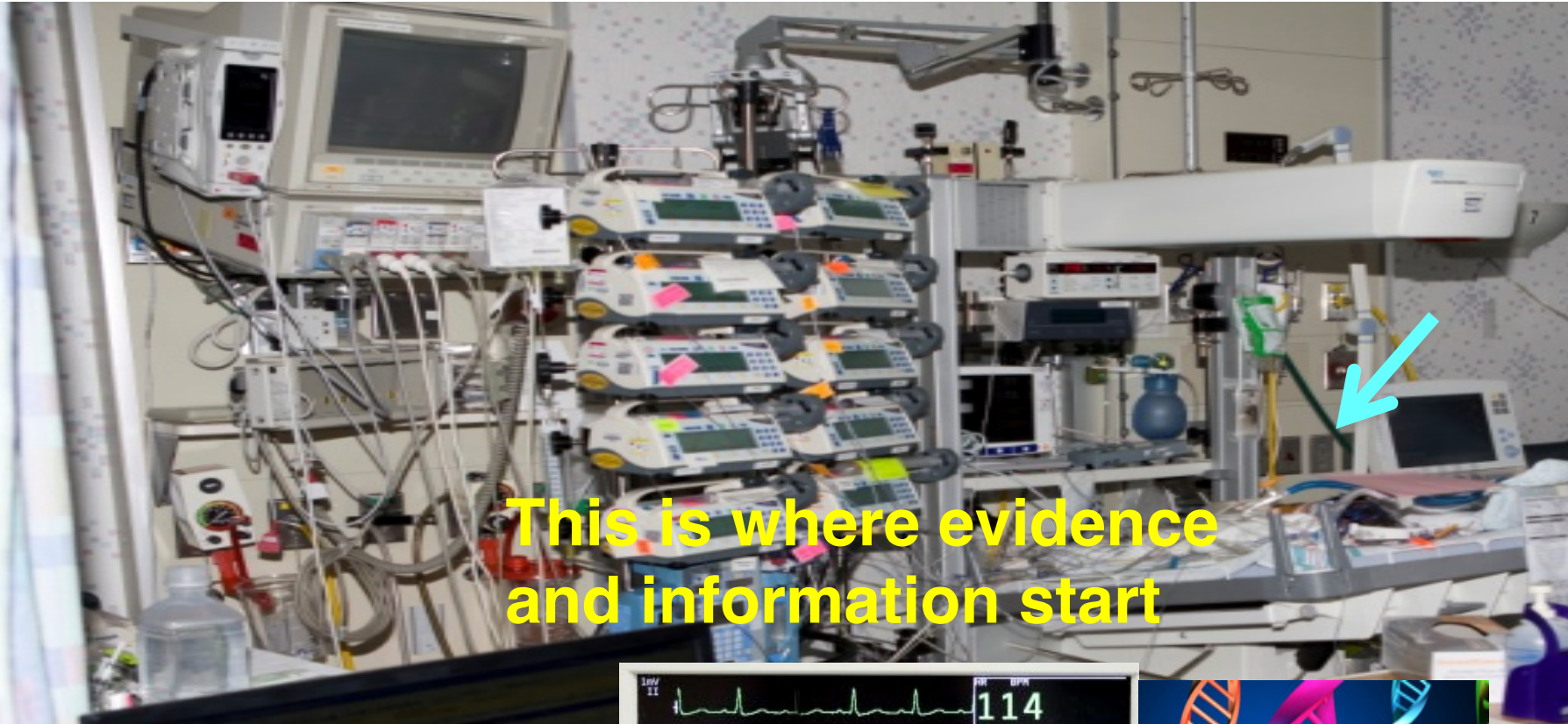
Dashboard



Factory IT Infrastructure




The Healthcare Industry



Building an ready-to-use AI solution for this is

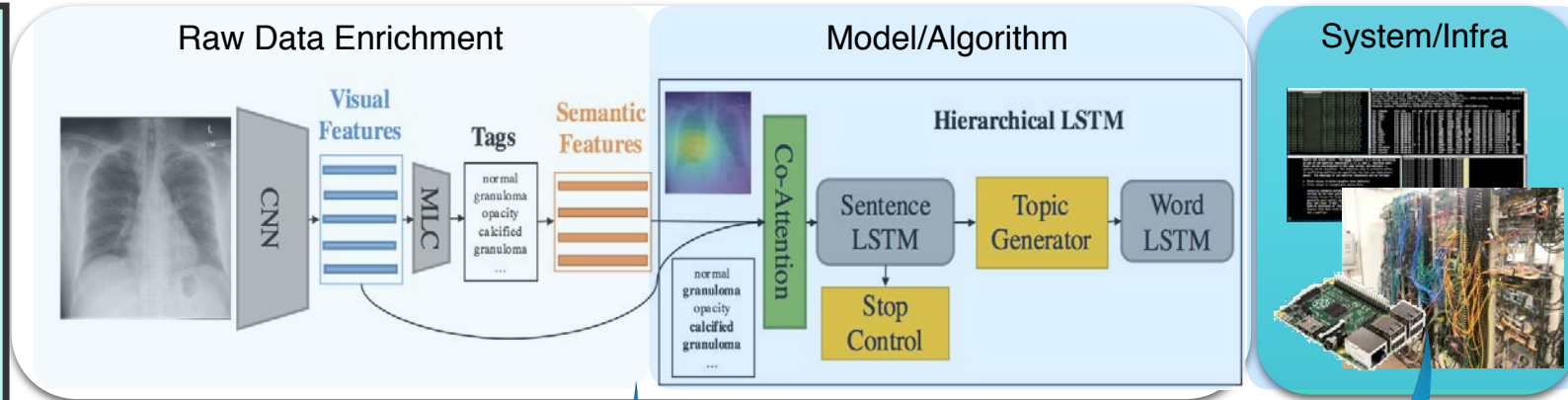
Extremely complex



Findings:
 There are no focal areas of consolidation.
 No suspicious pulmonary opacities.
 Heart size within normal limits.
 No pleural effusions.
 There is no evidence of pneumothorax.
 Degenerative changes of the thoracic spine.

Impression:
 No acute cardiopulmonary abnormality.

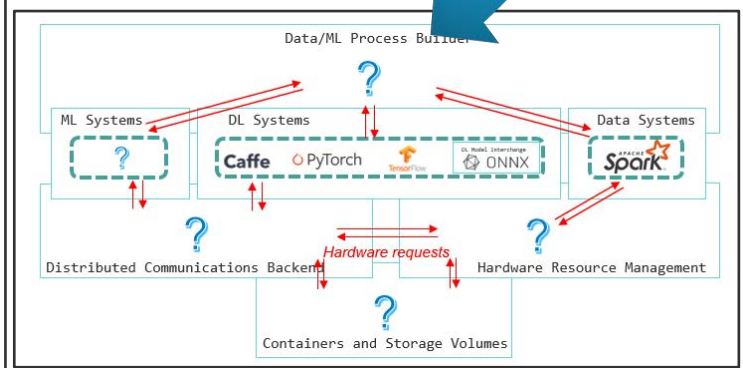
Task: Automatic Medical Report Generation



Requires inter-operation between diverse systems



User interface for Doctors



- GPU / CPU
- Mobile
- Laptop / PC
- Data Center
- IoT
- And More ...

AI as of now: Not Built, but Crafted



Outline

- Composable ML

A first-principle view of ML components and assemblage



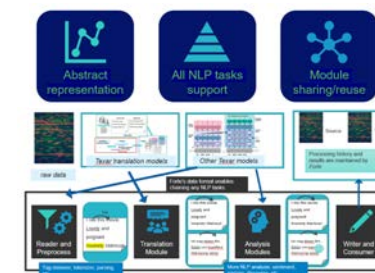
- Texar: A Modularized ML toolkit

Compose your ML applications like playing building blocks



- Scalable AI Infrastructure

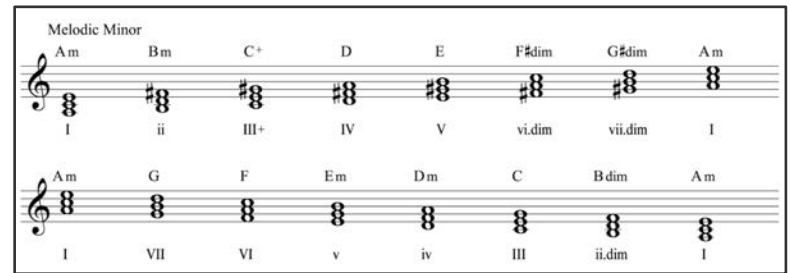
Composable ML in production



Composable ML



One-off design and programming

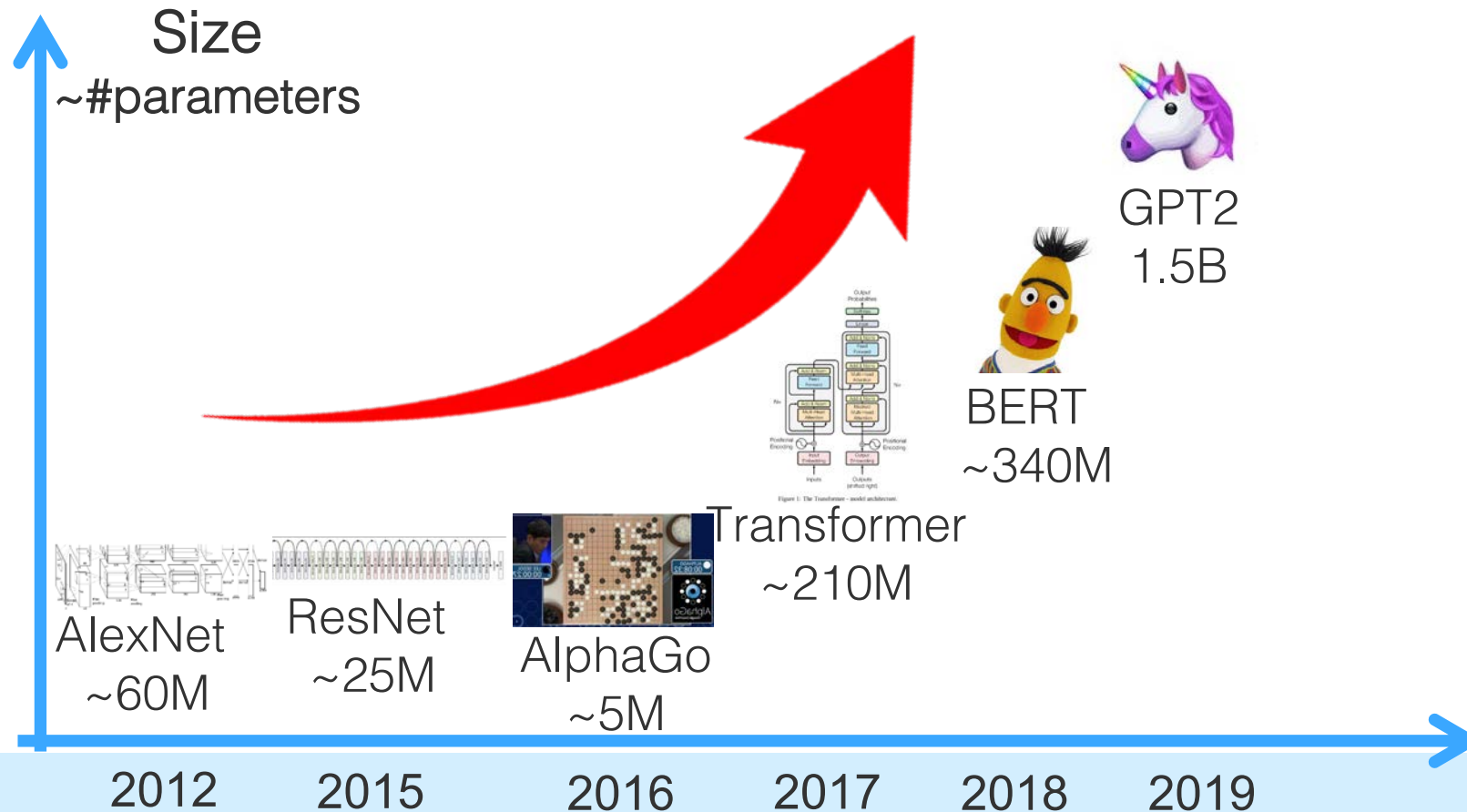


Modular and standardized

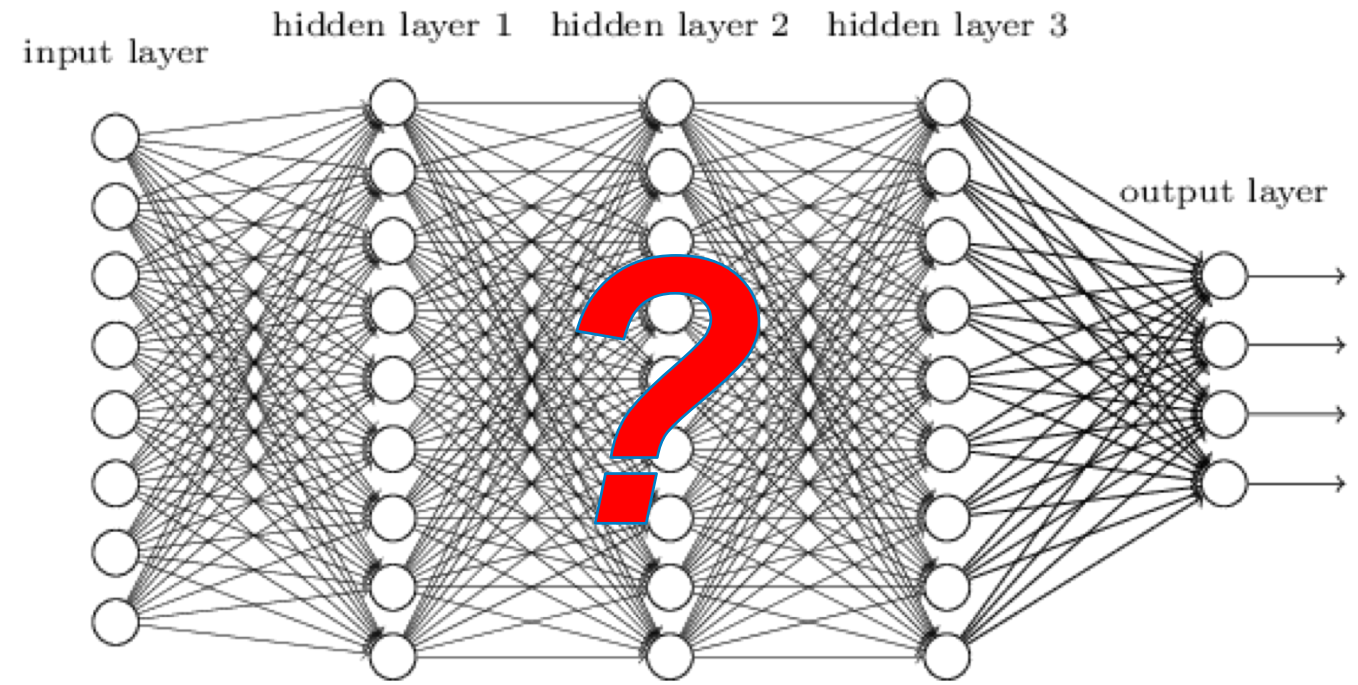
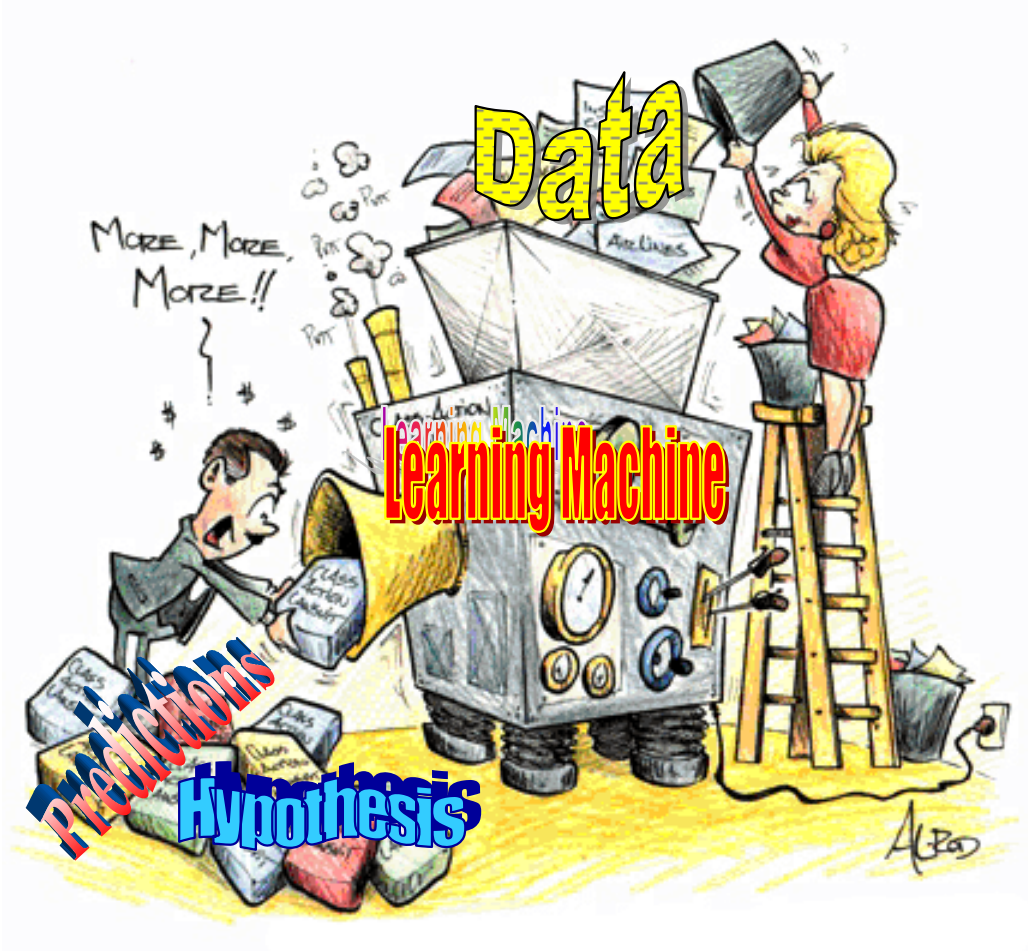
Complex “rhythms” and “chords”
systematically built out of simpler “notes”

Single Models with Increasing Size and Performance

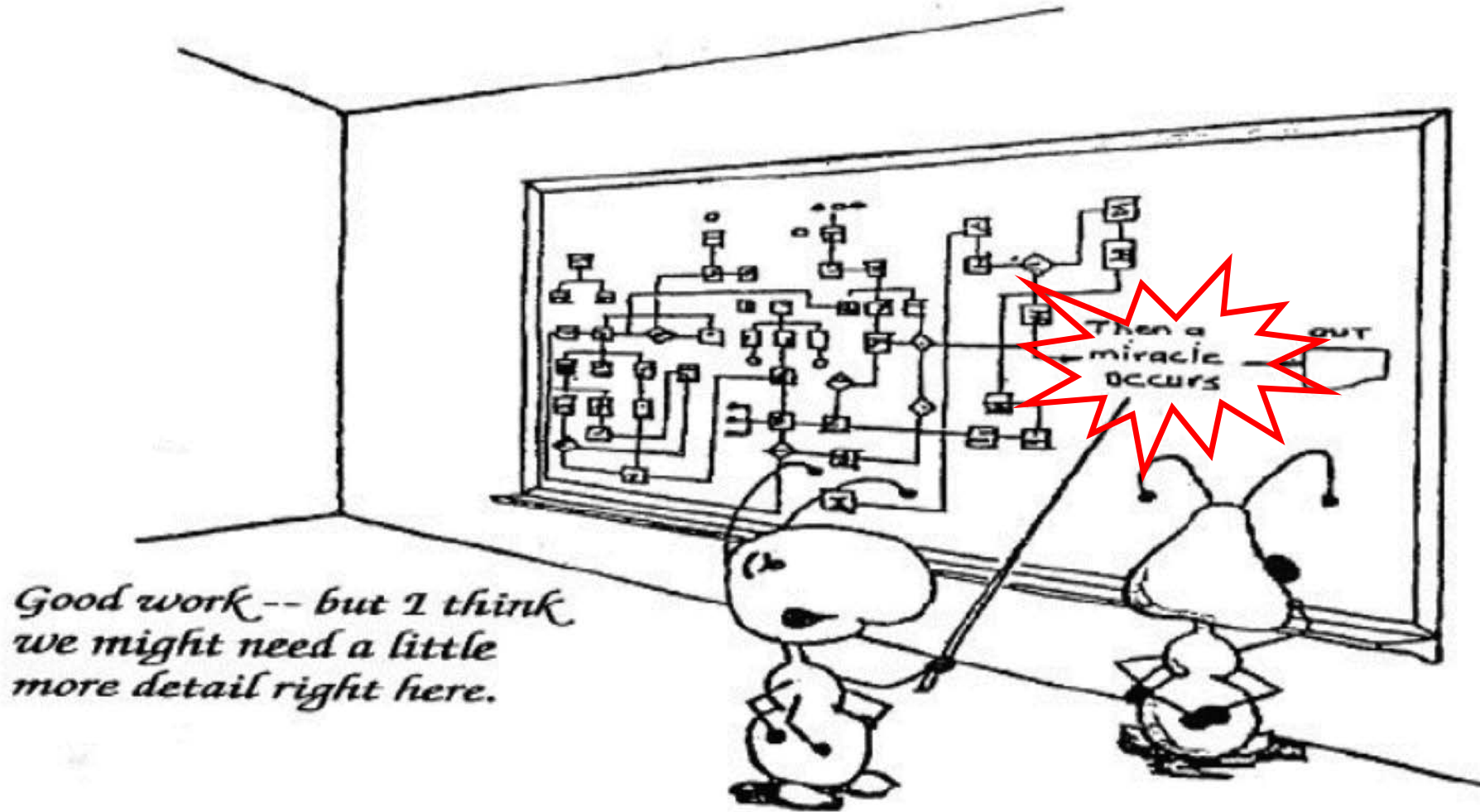
- Increasingly large black-box neural networks
- Good, even super-human performance on some tasks



Single Giant Models Enough?

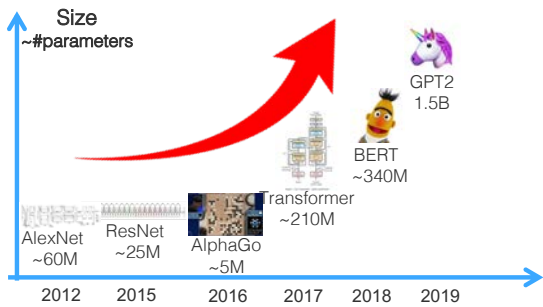


Difficulties of Single Giant Models



- Explainability
- Expertise
- Debugging
- Maintenance
- Upgrade
- Scalability
- ...

Far from Solving Real Complex Problems

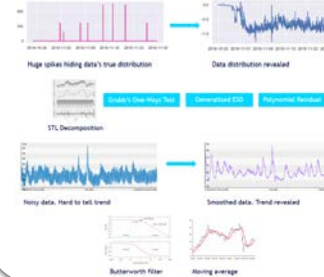


V.S.

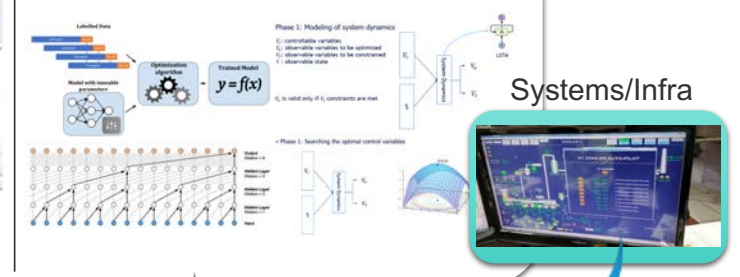
Extremely complex



Raw data far from ideal



Simultaneous Prediction and Control

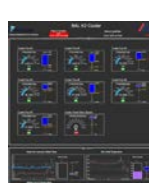


Requires inter-operation between diverse systems

Can't solve with "algo marketplace" or Kaggle competition

Extremely complex

Dashboard

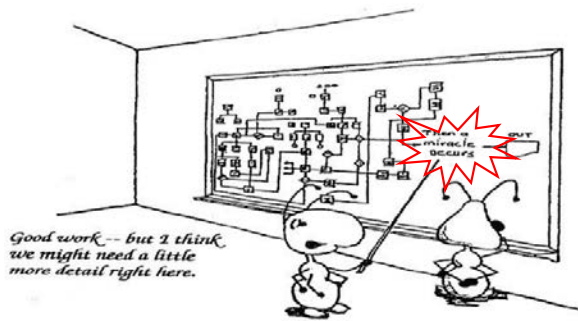
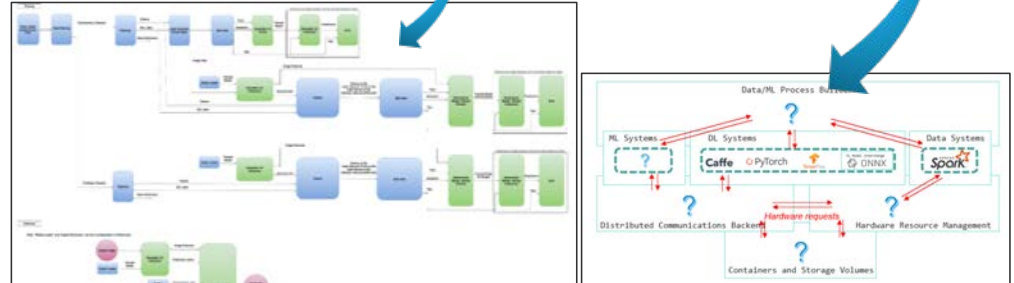
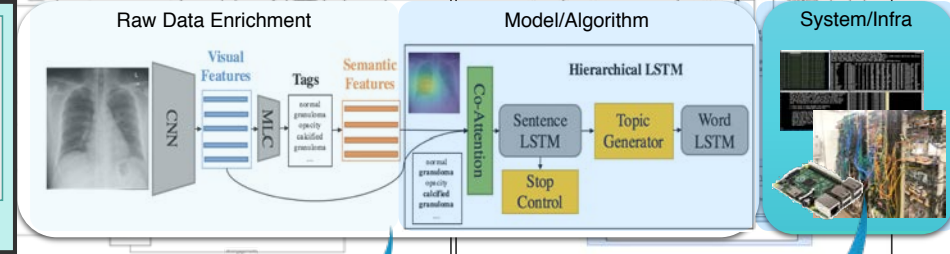


Findings:
There are no focal areas of consolidation.
No suspicious pulmonary opacities.
Heart size within normal limits.
No pleural effusions.
There is no evidence of pneumothorax.
Degenerative changes of the thoracic spine.
Impression:
No acute cardiopulmonary abnormality.

Task: Automatic Medical Report Generation

Requires inter-operation between diverse systems

User interface for Doctors



Good work-- but I think we might need a little more detail right here.

We choose MT as our running example to keep this talk to 90min

Running Example: Machine Translation (MT)



raw data



cleaning
tokenizing
vocabulary
truncation
...

source.dat
I like this movie.
Lovely and poignant
Insanely hilarious!
...

target.dat
Ich mag diesen film.
Schön und ergreifend
Wahnsinnig witzig!
...

clean data

training

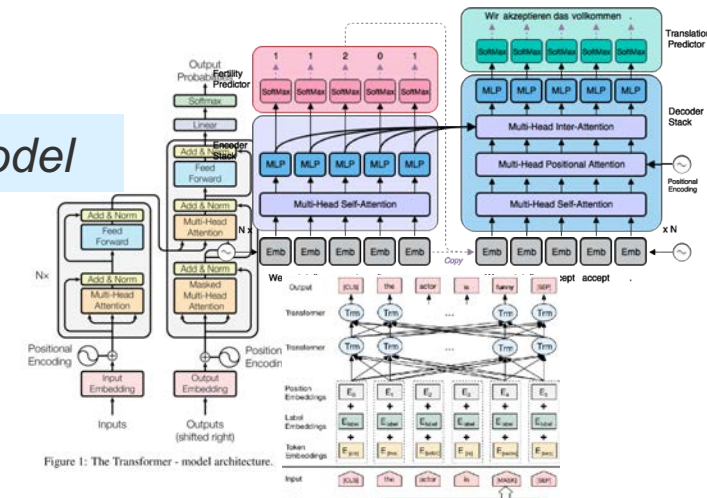
Maximum likelihood training

Reinforcement learning

Adversarial learning

Finetuning

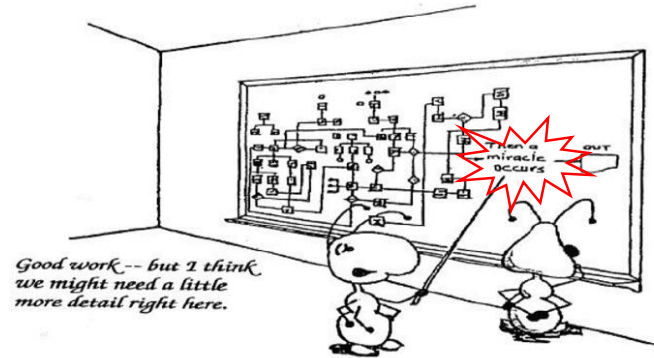
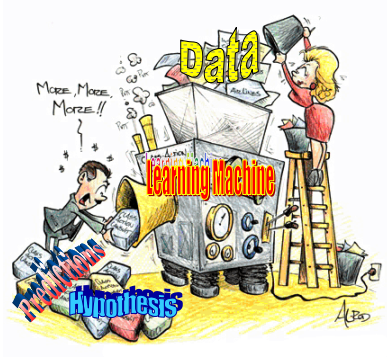
model



evaluation
post-processing
...



Solution: Composable ML



Composable ML



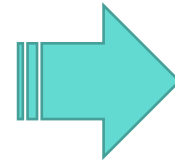
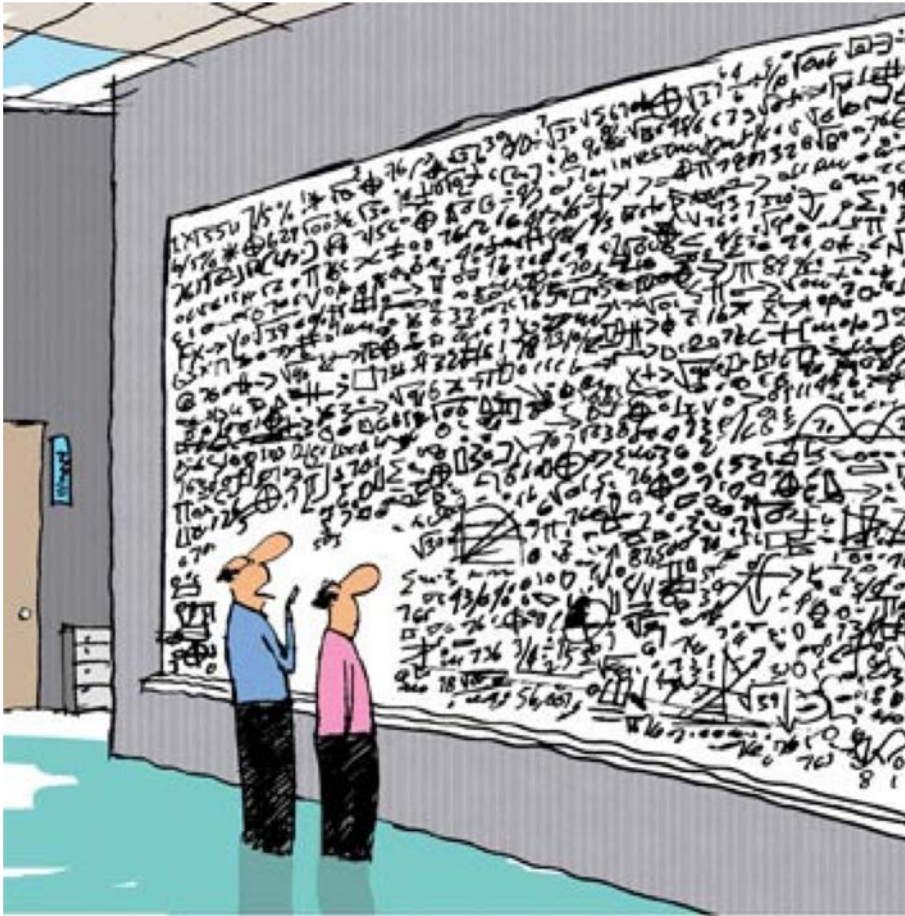
A modularized way to build complex applications

Solution: Composable ML

- Build AI solutions more easily, via pick-and-choose:
 - Data: text, speech, image, video, time series, ...
 - Models: recurrent, transformer, convolutional, ...
 - Tasks: classification, regression, generation, discovery, ...
- Stop writing same one-off code again and again
 - More reliable and easier to debug
 - Easier to onboard new developers



First Principles: Decomposing Machine Learning



Loss functions

(likelihood, reconstruction, margin, ...)

Model architectures

(RNNs, Transformers, Graphical, ...)

Constraints

(normality, sparsity, logical, KL, sum, ...)

Algorithms

MC (MCMC, Importance), Opt (gradient, IP), ...

Data

(processing, augmentation, weighting, ...)

Decomposing Machine Learning

Machine Learning:

Computational methods that enable machines to learn concepts and improve performance from experience

$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$



Decomposing Machine Learning

Machine Learning:

Computational methods that enable machines to learn concepts and improve performance from experience

$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$



model architecture/
inference procedure

$$y \sim p_{\theta}(y|x)$$

Decomposing Machine Learning

Machine Learning:

Computational methods that enable machines to learn concepts and improve performance from experience

$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

(x^*, y^*)

model architecture/
inference procedure

experience (data)

The diagram illustrates the decomposition of the machine learning optimization problem. The equation $\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$ is shown. A dashed arrow points from θ to the text 'model architecture/inference procedure'. Another dashed arrow points from D to the text 'experience (data)'. The optimal solution (x^*, y^*) is shown to the right of the equation.

Decomposing Machine Learning

Machine Learning:

Computational methods that enable machines to learn concepts and improve performance from experience

$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

loss

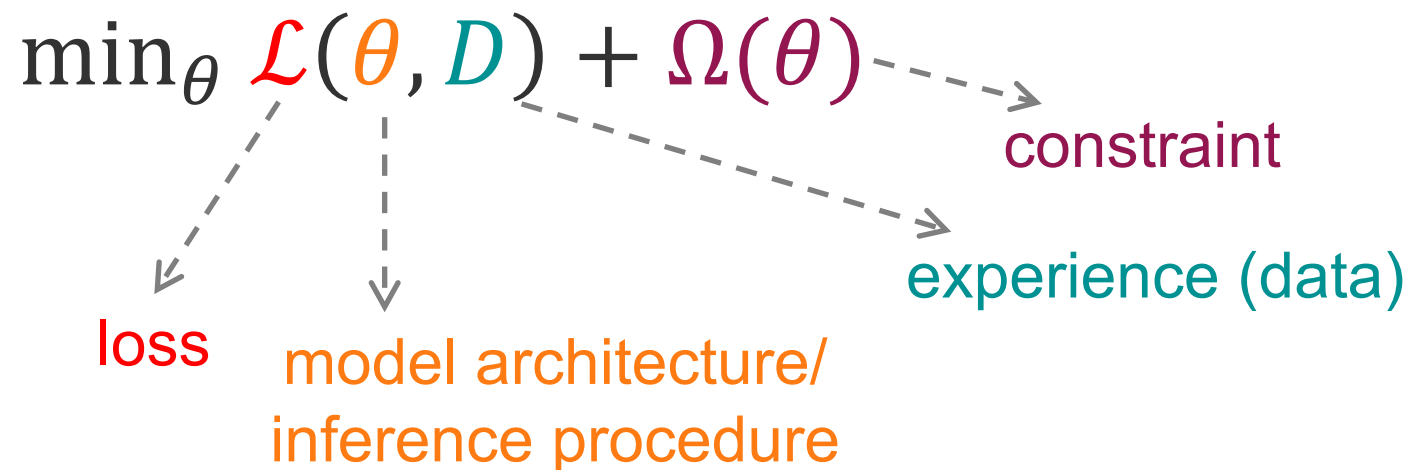
model architecture/
inference procedure

experience (data)

Decomposing Machine Learning

Machine Learning:

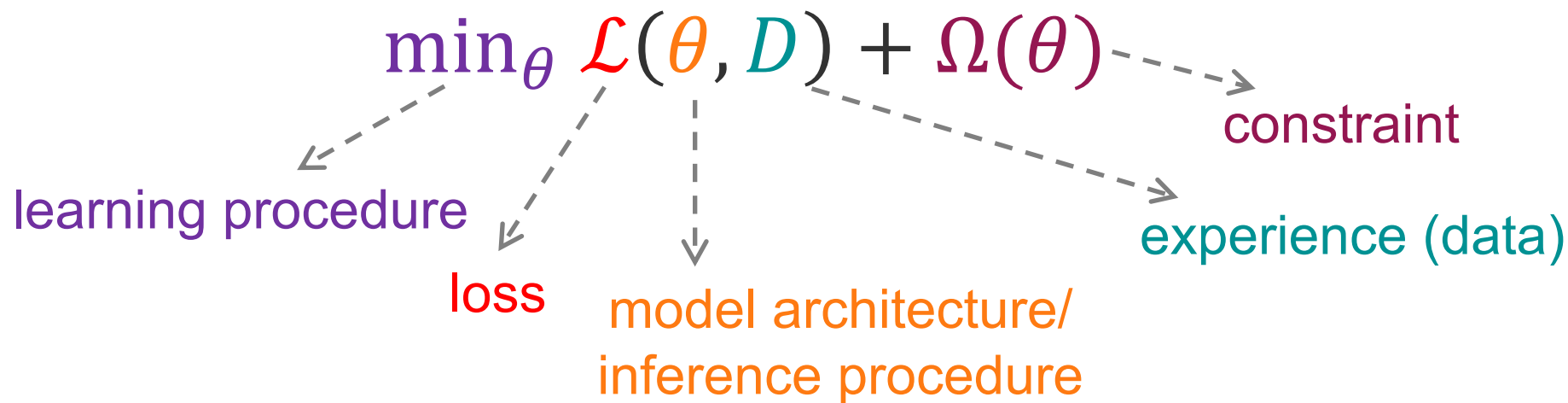
Computational methods that enable machines to learn concepts and improve performance from experience



Decomposing Machine Learning

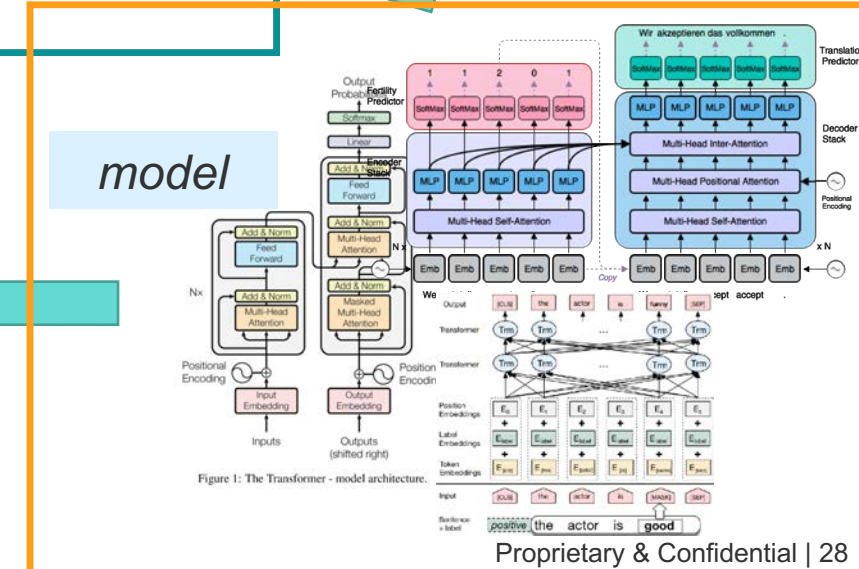
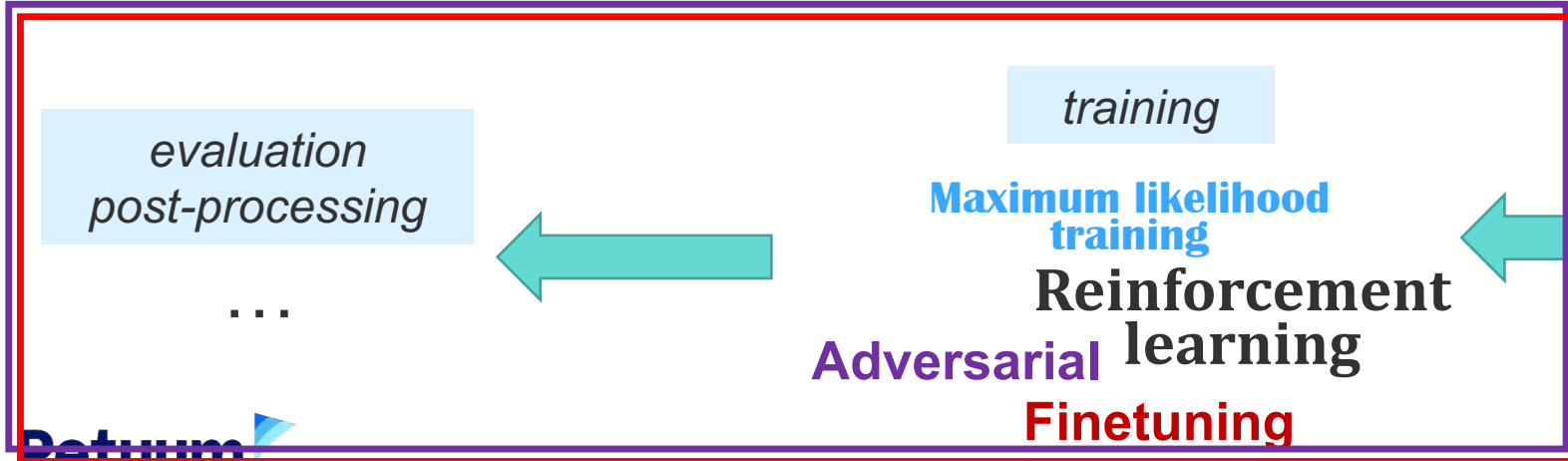
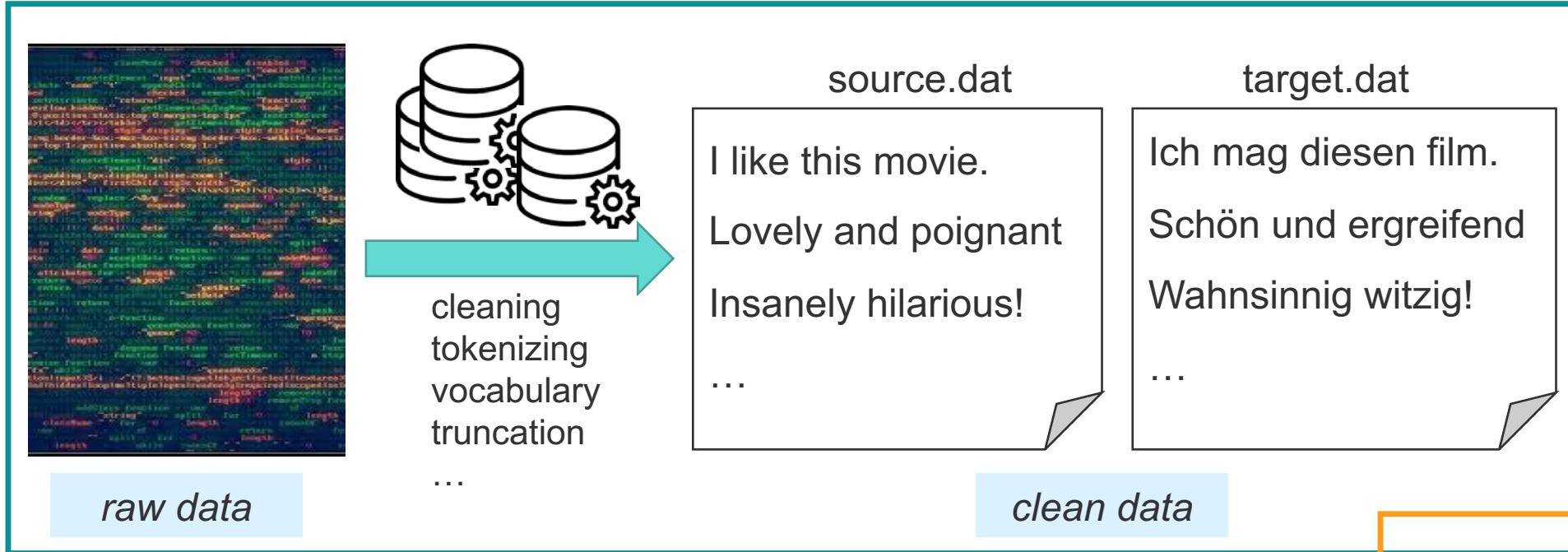
Machine Learning:

Computational methods that enable machines to learn concepts and improve performance from experience



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Running Example: Machine Translation



ML Components



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Constraint

Loss

Learning

Inference

Architecture

ML Components



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Constraint

Loss

Learning

Inference

Architecture

Architecture (1): Language Model



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

- Calculates the probability of a sentence:
 - Sentence:

$$\mathbf{y} = (y_1, y_2, \dots, y_T)$$

Example:

(I, like, this, ...)

Architecture (1): Language Model



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

- Calculates the probability of a sentence:
 - Sentence:

$$\mathbf{y} = (y_1, y_2, \dots, y_T)$$

$$p_{\theta}(\mathbf{y}) = \prod_{t=1}^T p_{\theta}(y_t | \mathbf{y}_{1:t-1})$$

Example:

(I, like, this, ...)

$\dots p_{\theta}(\textit{like} | I) p_{\theta}(\textit{this} | I, \textit{like}) \dots$



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Architecture (1): Language Model

- Calculates the probability of a sentence:
 - Sentence:

$$\mathbf{y} = (y_1, y_2, \dots, y_T)$$

$$p_{\theta}(\mathbf{y}) = \prod_{t=1}^T p_{\theta}(y_t | \mathbf{y}_{1:t-1})$$

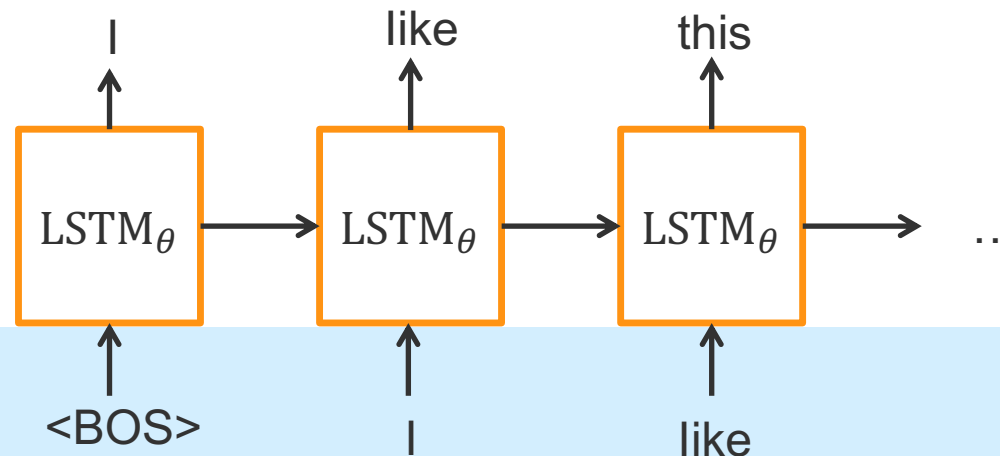
Example:

(I, like, this, ...)

$\dots p_{\theta}(\text{like} | I) p_{\theta}(\text{this} | I, \text{like}) \dots$

Architecture (1.1)

LSTM RNN





$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Architecture (1): Language Model

- Calculates the probability of a sentence:
 - Sentence:

$$\mathbf{y} = (y_1, y_2, \dots, y_T)$$

$$p_{\theta}(\mathbf{y}) = \prod_{t=1}^T p_{\theta}(y_t | \mathbf{y}_{1:t-1})$$

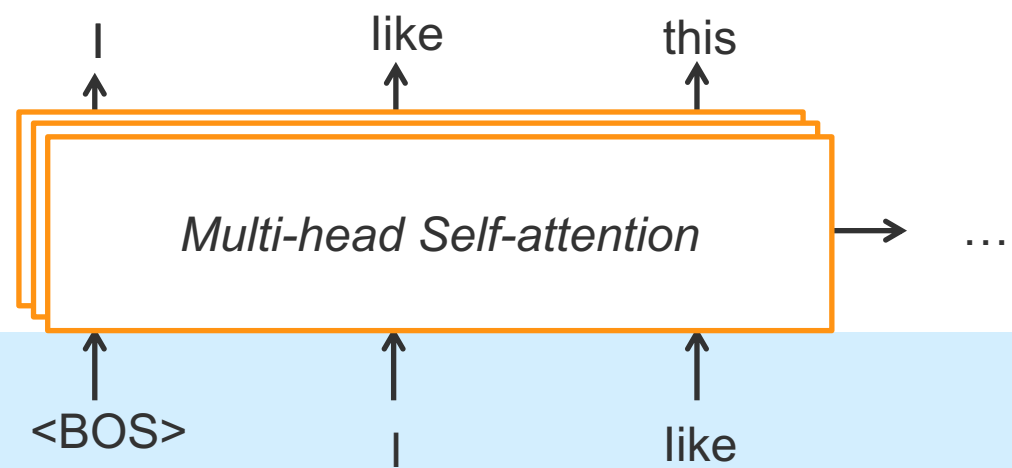
Example:

(I, like, this, ...)

$\dots p_{\theta}(\textit{like} | I) p_{\theta}(\textit{this} | I, \textit{like}) \dots$

Architecture (1.2)

Transformer



Architecture (2): **Conditional** Language Model



- Conditions on additional task-dependent context x
 - Machine translation: source sentence

I like this movie. → Ich mag diesen film.

- Medical image report generation: medical image



... There is chronic pleural-parenchymal scarring within the lung bases. No lobar consolidation is seen. ...

Architecture (2): **Conditional** Language Model



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

- Conditions on additional task-dependent context \mathbf{x}

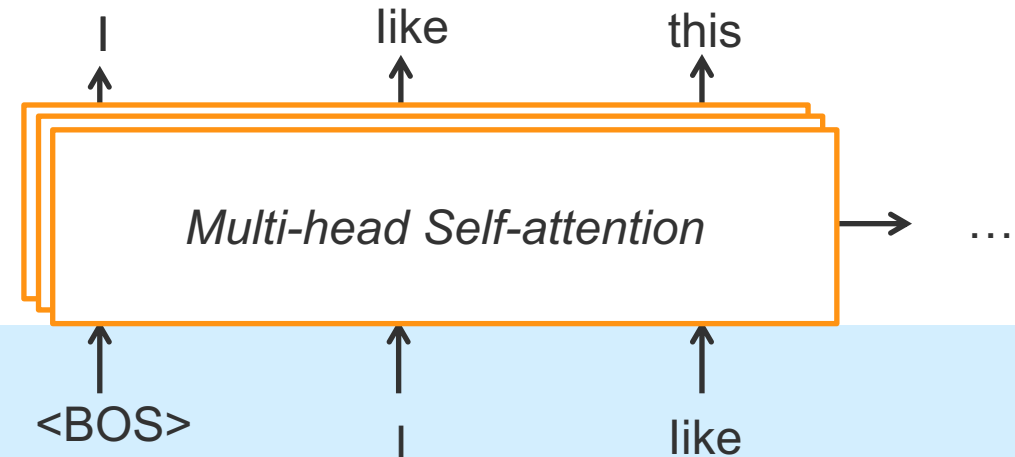
$$p_{\theta}(\mathbf{y} | \mathbf{x}) = \prod_{t=1}^T p_{\theta}(y_t | \mathbf{y}_{1:t-1}, \mathbf{x})$$

Architecture (2): **Conditional** Language Model



- Conditions on additional task-dependent context \mathbf{x}

$$p_{\theta}(\mathbf{y} | \mathbf{x}) = \prod_{t=1}^T p_{\theta}(y_t | \mathbf{y}_{1:t-1}, \mathbf{x})$$



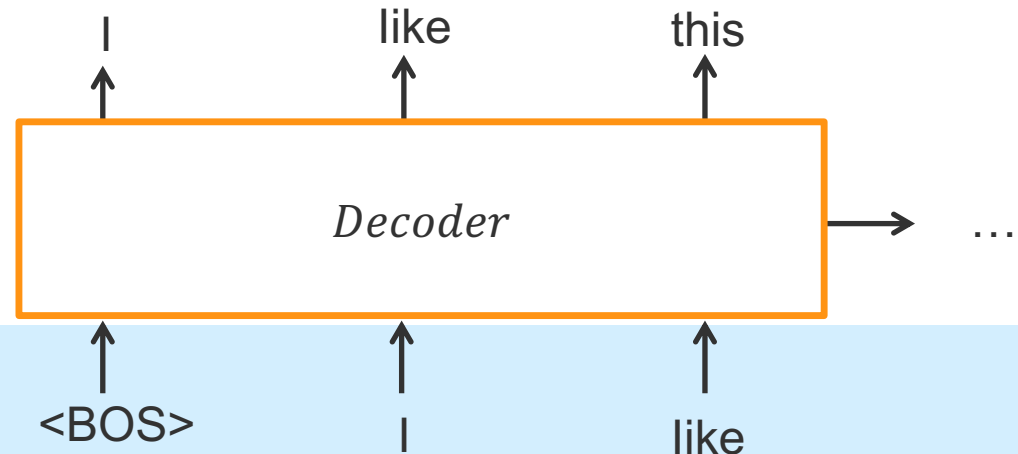
Architecture (2): **Conditional** Language Model



- Conditions on additional task-dependent context \mathbf{x}

$$p_{\theta}(\mathbf{y} | \mathbf{x}) = \prod_{t=1}^T p_{\theta}(y_t | \mathbf{y}_{1:t-1}, \mathbf{x})$$

- Language model as a **decoder**



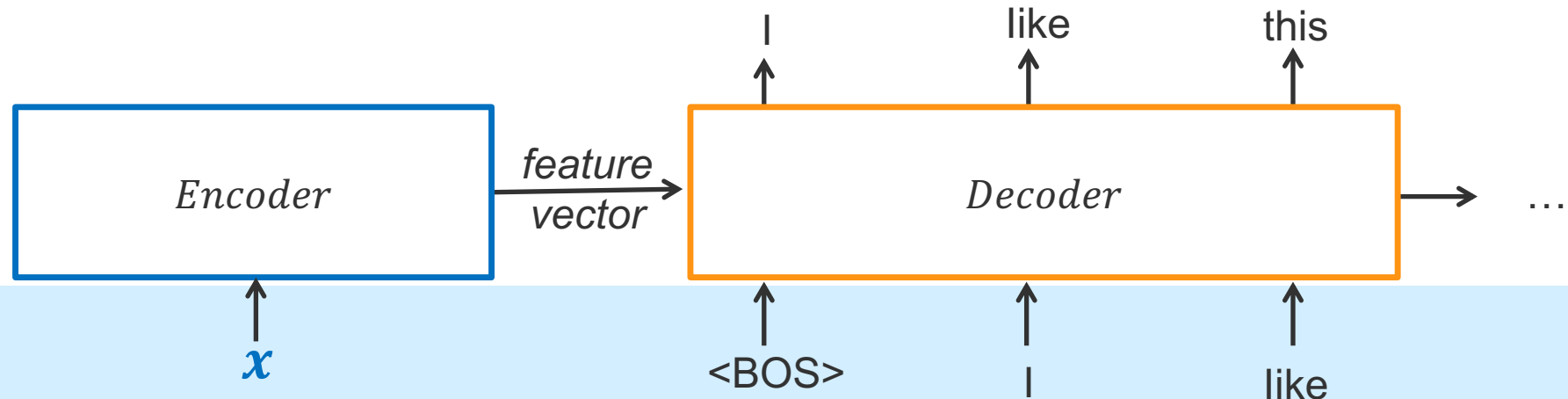
Architecture (2): **Conditional** Language Model



- Conditions on additional task-dependent context \mathbf{x}

$$p_{\theta}(\mathbf{y} | \mathbf{x}) = \prod_{t=1}^T p_{\theta}(y_t | \mathbf{y}_{1:t-1}, \mathbf{x})$$

- Language model as a **decoder**
- Encodes context with an **encoder**



Architecture Graph



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

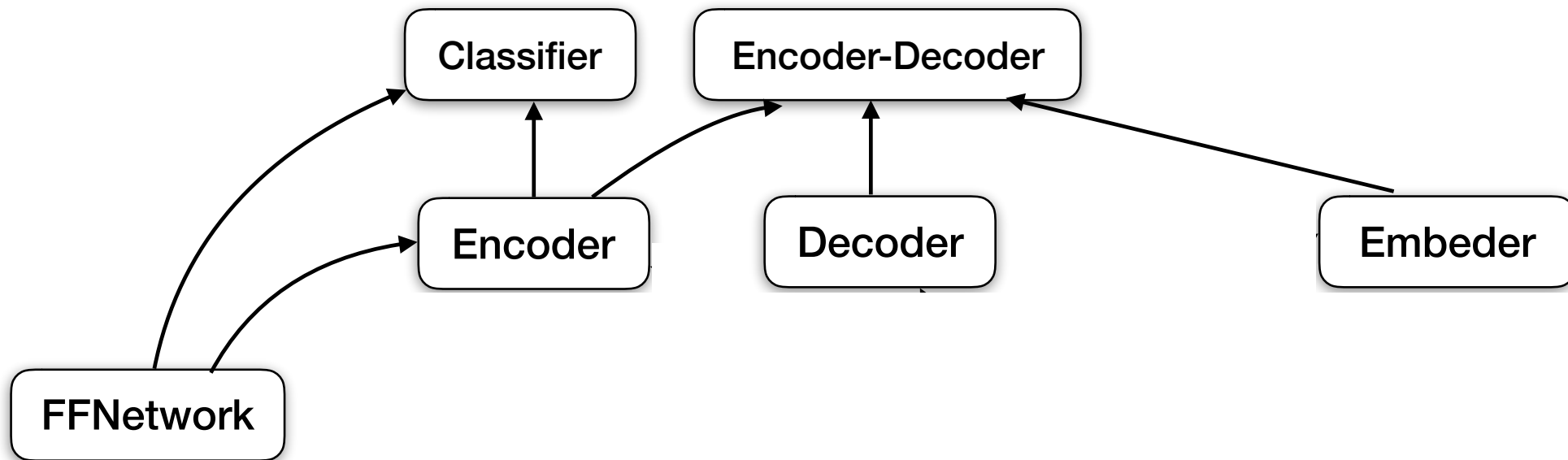
Classifier

Encoder-Decoder

Architecture Graph



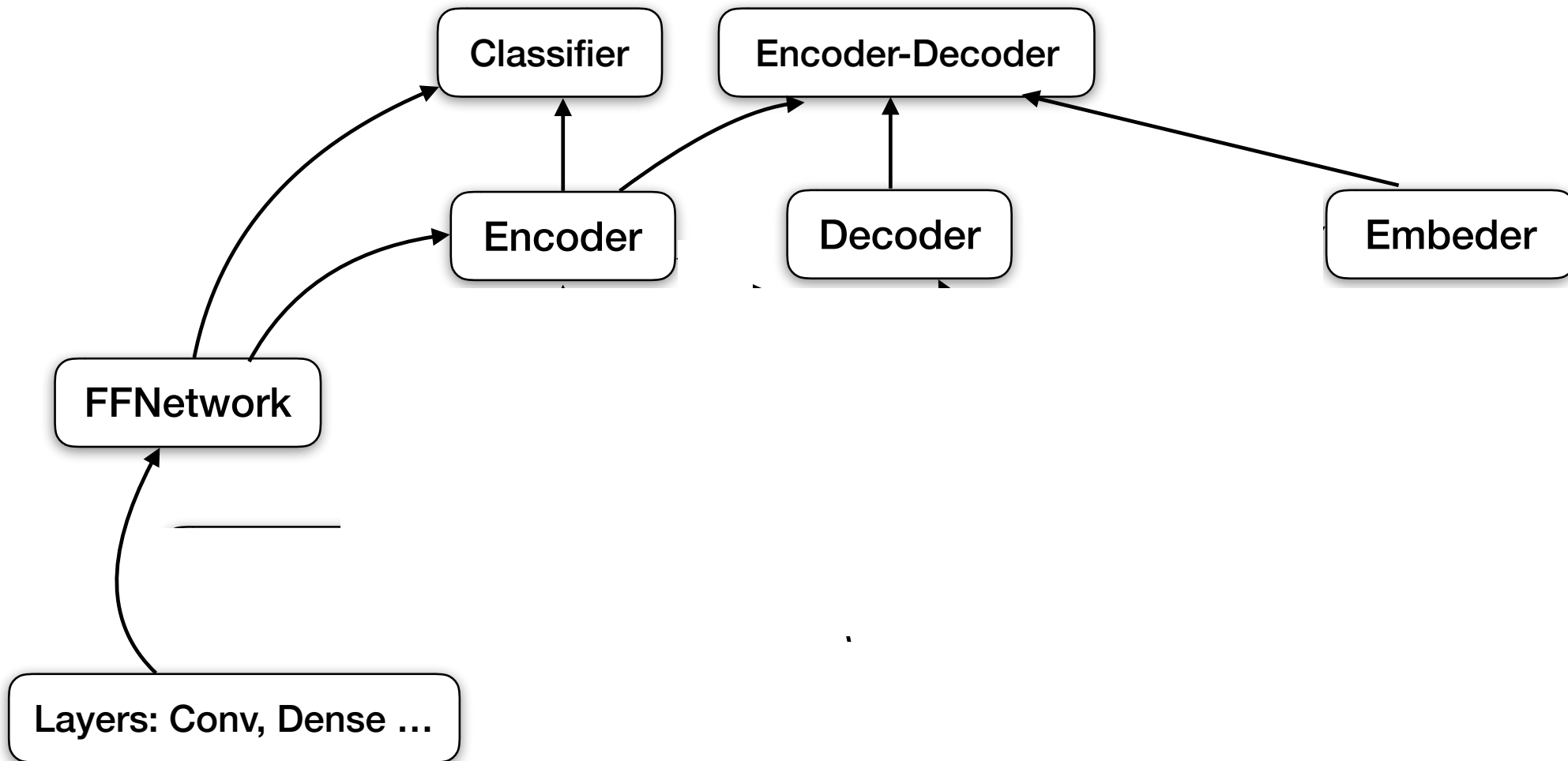
$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$



Architecture Graph



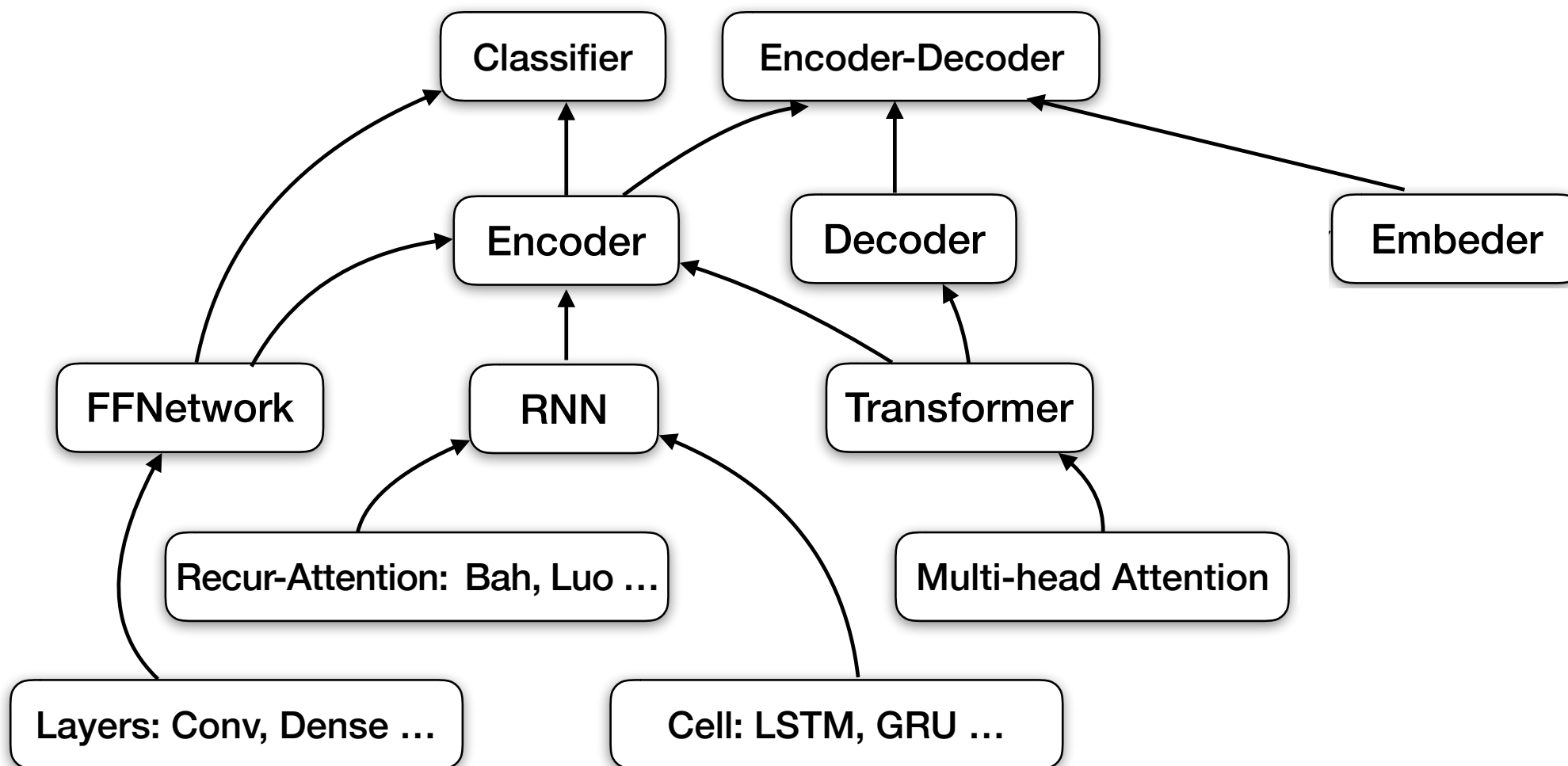
$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$



Architecture Graph



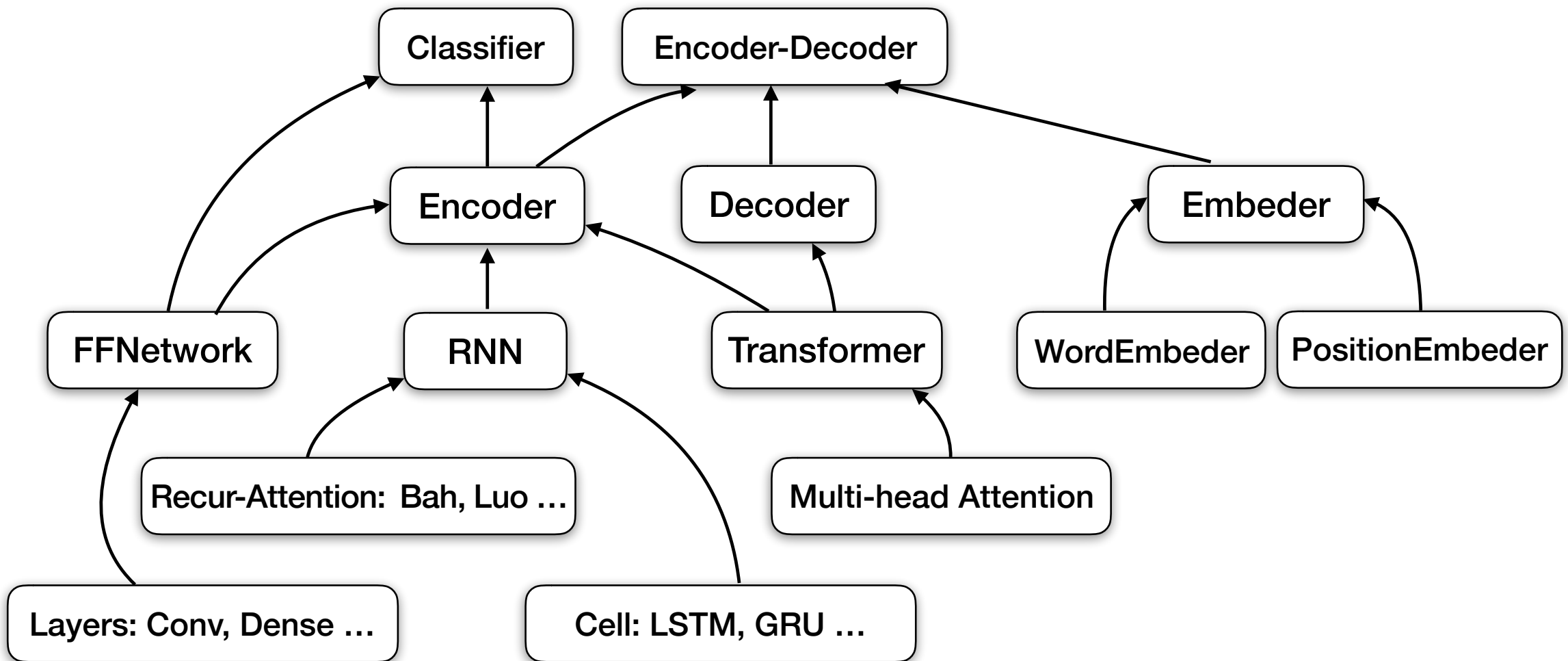
$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$



Architecture Graph



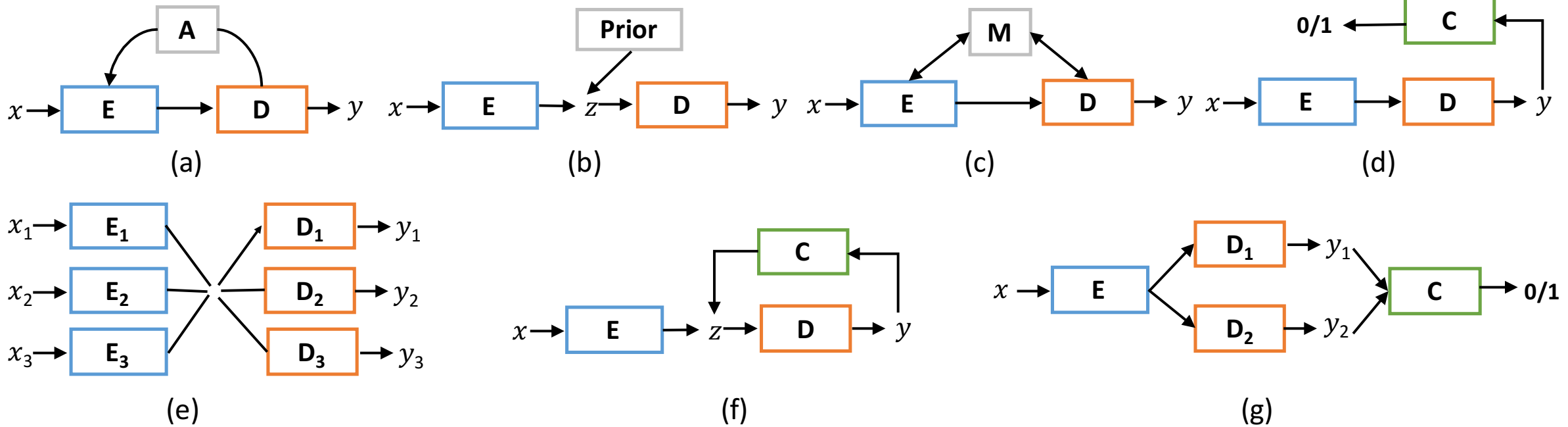
$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$



Complex Composite Architectures



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$



E refers to encoder, D to decoder, C to Classifier, A to attention, Prior to prior distribution, and M to memory

ML Components



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Constraint

Loss

Learning

Inference

Architecture

decoder

LSTM RNN

Attention RNN

Transformer

...

encoder

classifier

...

ML Components



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Constraint

Loss

Learning

Inference

Architecture

decoder

LSTM RNN

Attention RNN

Transformer

...

encoder

classifier

...

Learning, Inference & Loss (1): Maximum Likelihood Estimation



- Given data example $(\mathbf{x}^*, \mathbf{y}^*)$
- Maximizes log-likelihood of the data

Learning

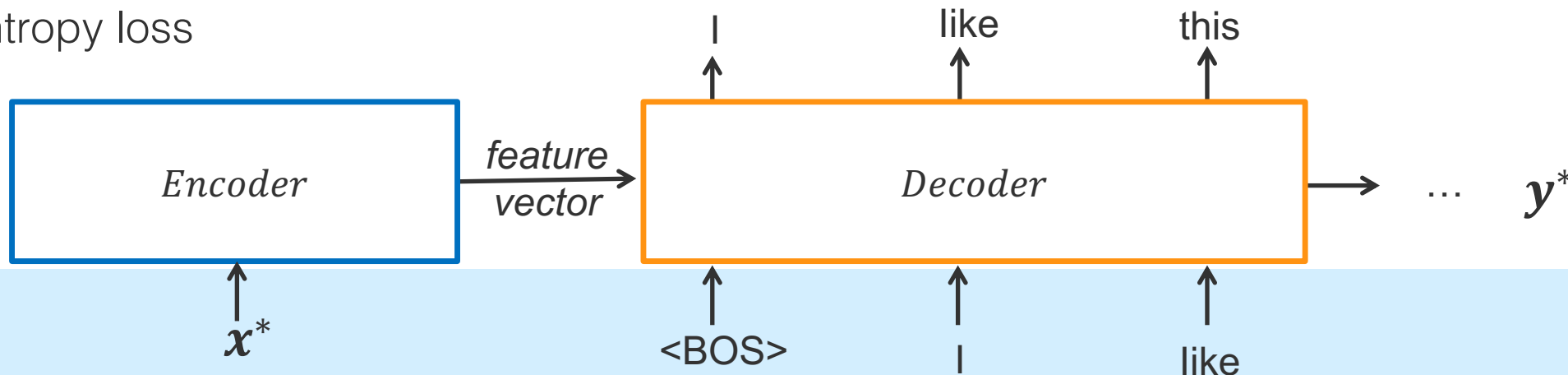
$$\begin{aligned} \min_{\theta} \mathcal{L}_{\text{MLE}} &= -\log p_{\theta}(\mathbf{y}^* | \mathbf{x}^*) \\ &= -\prod_{t=1}^T p_{\theta}(y_t^* | \mathbf{y}_{1:t-1}^*, \mathbf{x}^*) \end{aligned}$$

Loss Cross-entropy loss

Inference

Teacher-forcing decoding:

For every step t , feeds in the previous ground-truth tokens $y_{1:t-1}^*$ to decode next step



Learning, Inference & Loss (2): Adversarial Learning



Learning

- A **discriminator** is trained to distinguish b/w *real data examples* and *fake generated samples*
- The **model** is trained to fool the discriminator

Loss

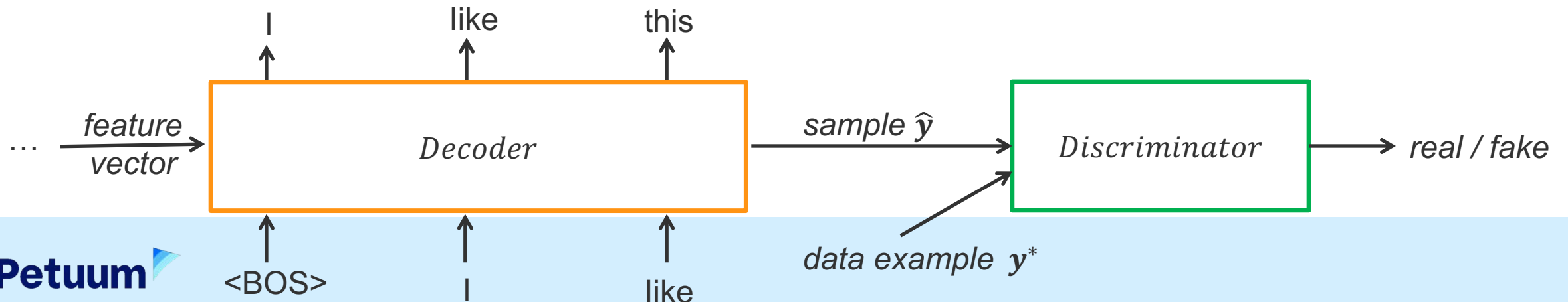
- Binary adversarial loss
- Feature-matching adversarial loss

Inference

Gumbel-softmax decoding:

Uses a differentiable approximation of sample \hat{y} for gradient backpropagation

$$\frac{\partial \mathcal{L}(\hat{y})}{\partial \theta} = \frac{\partial \mathcal{L}(\hat{y})}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \theta}$$



Learning, Inference & Loss (3): Reinforcement Learning



Learning

- Optimizes *expected* task reward
 - $\text{BLEU}(\hat{\mathbf{y}}, \mathbf{y}^*)$ for machine translation

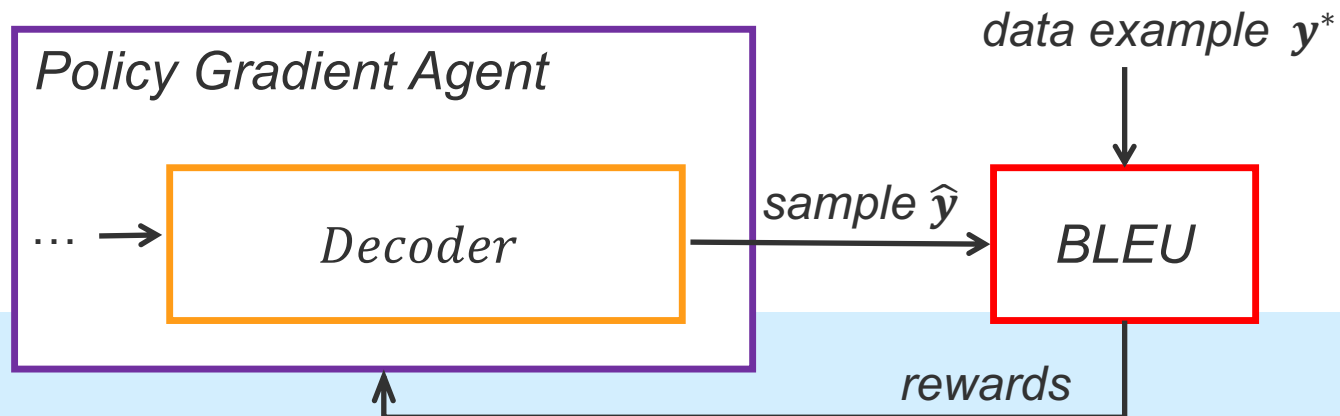
$$\mathbb{E}_{\hat{\mathbf{y}} \sim p_{\theta}(\mathbf{y} | \mathbf{x})} [\text{BLEU}(\hat{\mathbf{y}}, \mathbf{y}^*)]$$

Loss

- Policy gradient loss
- Policy gradient loss w/ baseline
- ...

Inference

- Greedy decoding
- Sampling decoding
- Beam search decoding
- Top- k / Top- p decoding
- ...



ML Components



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Constraint

Loss

Learning

Inference

Architecture

Cross-entropy

MLE

Teacher-forcing

decoder

Binary Adv loss

Adversarial

Gumbel-softmax

LSTM RNN

Matching Adv loss

Reinforcement

Sample

Attention RNN

PG loss

Adv + RL

Greedy

Transformer

PG loss + baseline

Structured super.

Beam-search

...

...

Reward-aug.

Top-k sample

encoder

...

...

classifier

...

ML Components



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Constraint

Loss

Learning

Inference

Architecture

Cross-entropy

MLE

Teacher-forcing

decoder

Binary Adv loss

Adversarial

Gumbel-softmax

LSTM RNN

Matching Adv loss

Reinforcement

Sample

Attention RNN

PG loss

Adv + RL

Greedy

Transformer

PG loss + baseline

Structured super.

Beam-search

...

...

Reward-aug.

Top-k sample

encoder

...

...

classifier

...

Constraint (1): Conventional Constraints

- Many choices for get different statistical properties:
 - Normality, Sparsity, KL, sum, ...



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Constraint (2): Structured Knowledge



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

- Structured knowledge as constraints

Sentiment classification:

- “Food was good, **but** the service was very disappointing.”

Logic rule: $f(x = \text{sentence}, y = \text{sentiment}) = \text{truth value}$

- Sentence x with structure $A\text{-but-}B \Rightarrow$ sentiment of B dominates



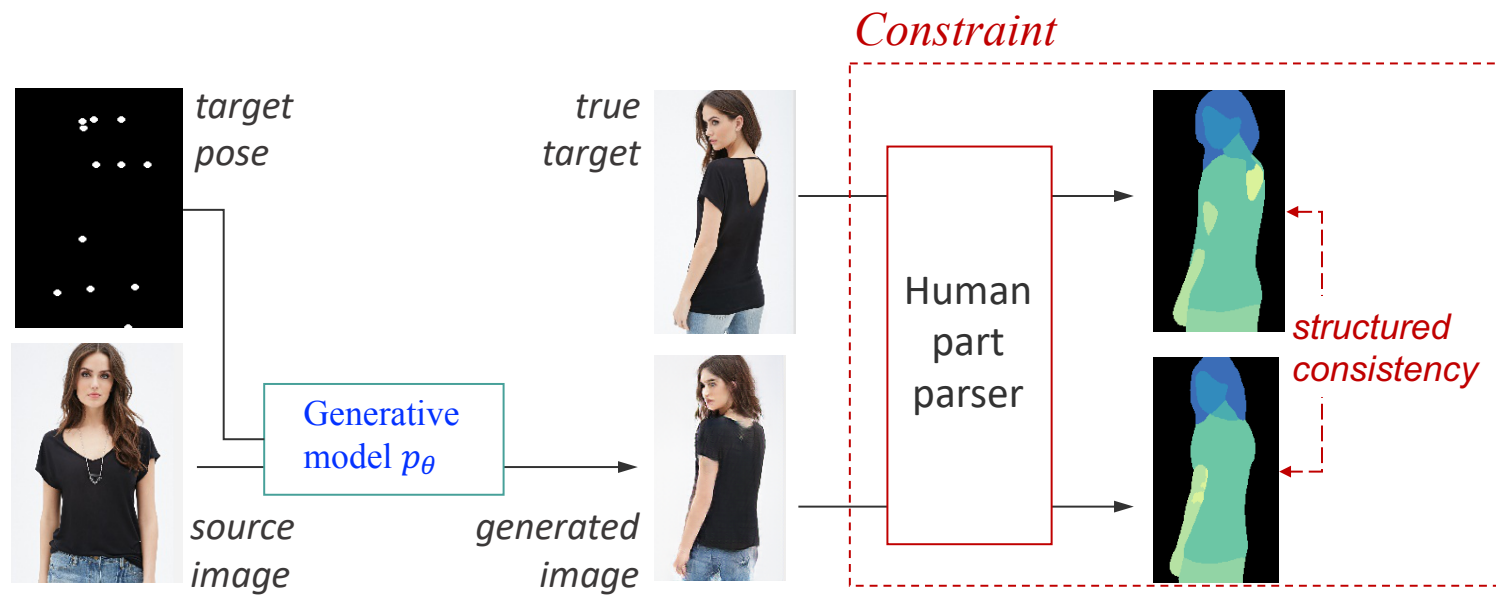
$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Constraint (2): Structured Knowledge

- Structured knowledge as constraints

Human Image Generation

$f(\mathbf{y} = \text{generated}, \mathbf{o} = \text{ground truth}) = \text{match score}$



Constraint (2): Structured Knowledge



- Structured knowledge as constraints

- Constraint function: $f(\mathbf{x}, \mathbf{y}, \mathbf{o}) \in \mathbb{R}$

- Model: $p_{\theta}(\mathbf{y}|\mathbf{x})$

- Variational Knowledge Regularization [Hu et al., 2016, 2018; Gal] Variational distribution

$$\mathcal{L}_{VKR}(q, \theta) = \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} [f(\mathbf{x}, \mathbf{y}, \mathbf{o})] - \text{KL}(q(\mathbf{y}|\mathbf{x}) || p_{\theta}(\mathbf{y}|\mathbf{x}))$$

- Learning procedure: at iteration n

$$q^{n+1}(\mathbf{y}|\mathbf{x}) \propto p_{\theta^n}(\mathbf{y}|\mathbf{x}) \exp\{f(\mathbf{x}, \mathbf{y}, \mathbf{o})\}$$

Combines the model and the knowledge -- teacher model

$$\theta^{n+1} = \operatorname{argmax}_{\theta} \mathbb{E}_{q^{n+1}(\mathbf{y}|\mathbf{x})} [\log p_{\theta}(\mathbf{y}|\mathbf{x})]$$

The model imitates the teacher model predictions -- student model

ML Components



$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

Constraint

L1 / L2

Logical

Structured

...

Loss

Cross-entropy

Binary Adv loss

Matching Adv loss

PG loss

PG loss + baseline

...

Learning

MLE

Adversarial

Reinforcement

Adv + RL

Structured super.

Reward-aug.

Inference

Teacher-forcing

Gumbel-softmax

Sample

Greedy

Beam-search

Top-k sample

...

Architecture

decoder

LSTM RNN

Attention RNN

Transformer

...

encoder

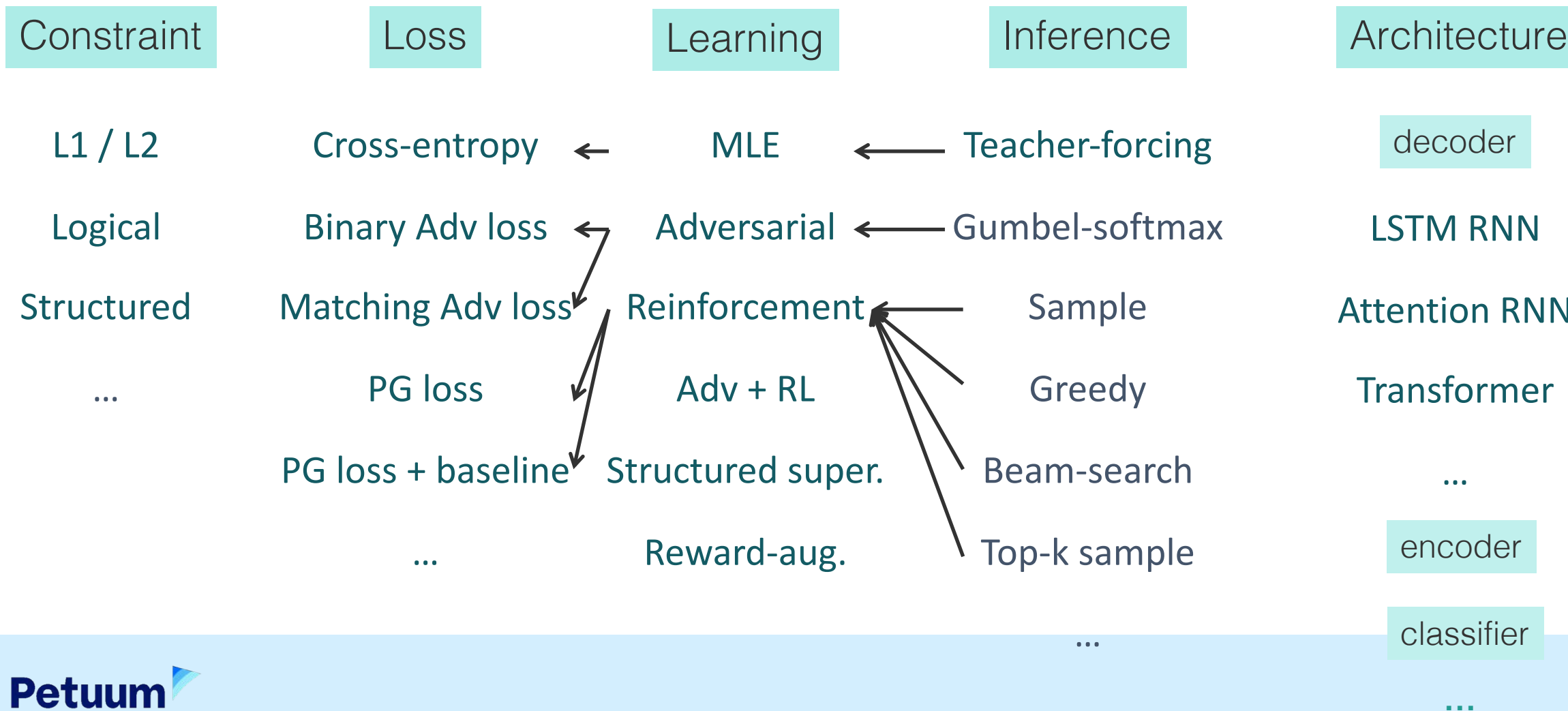
classifier

...

Holistic View



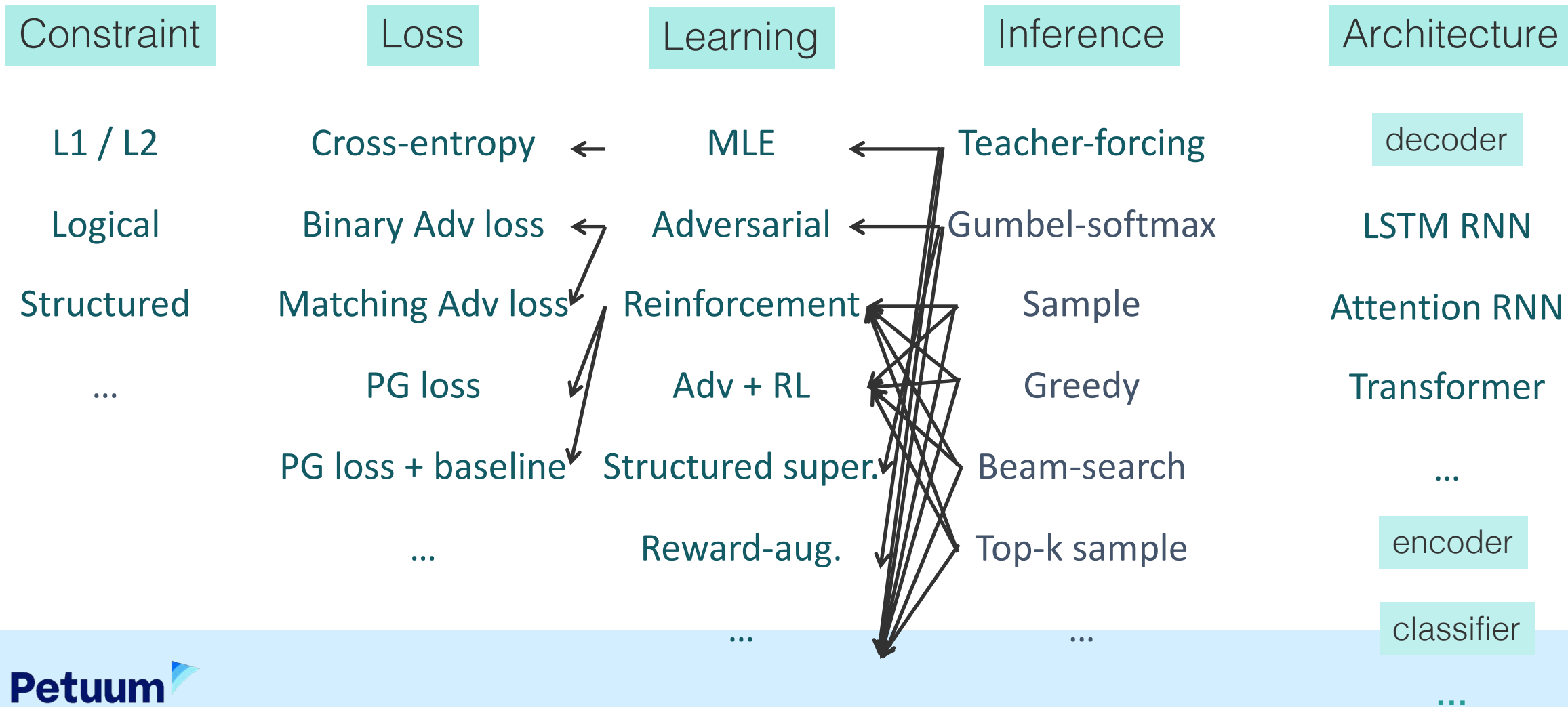
$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$



Holistic View



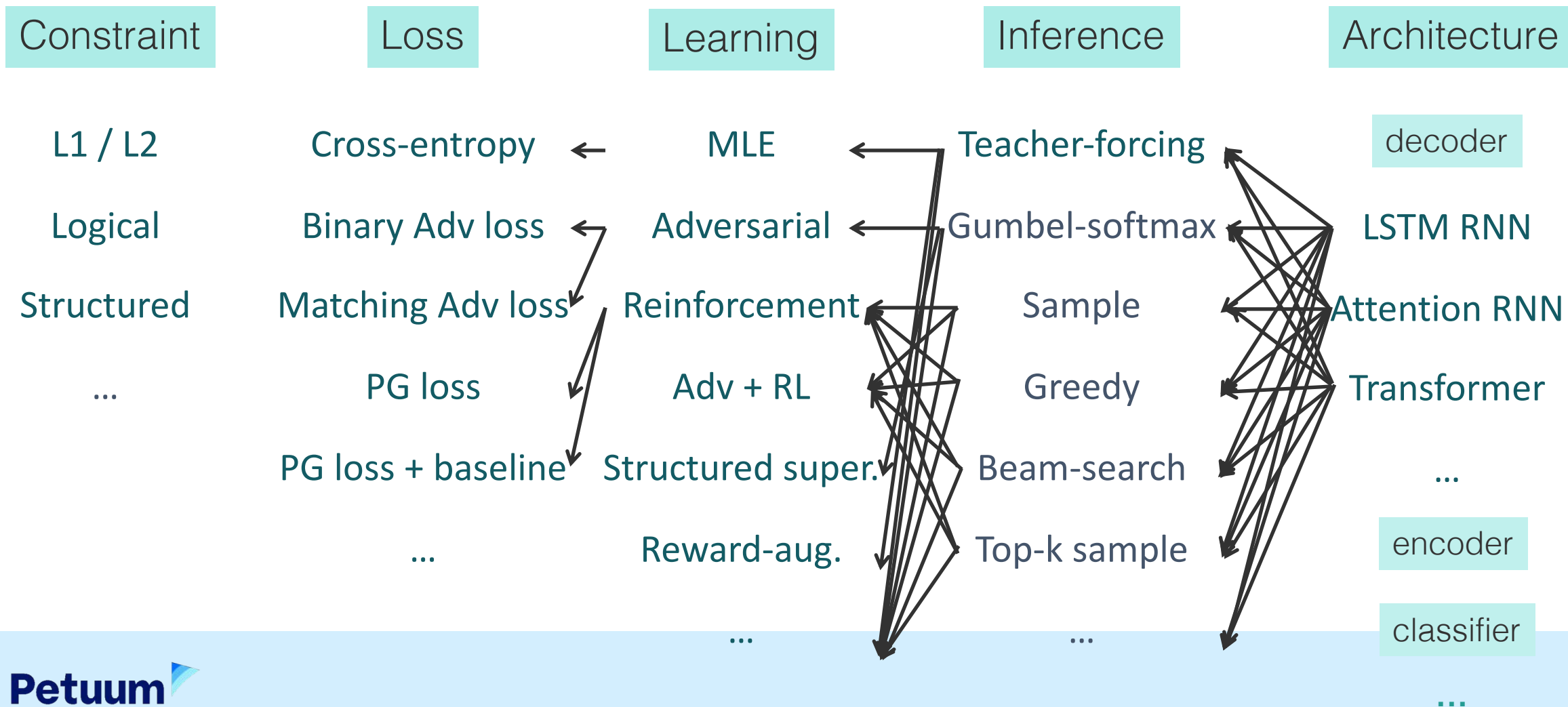
$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$



Holistic View



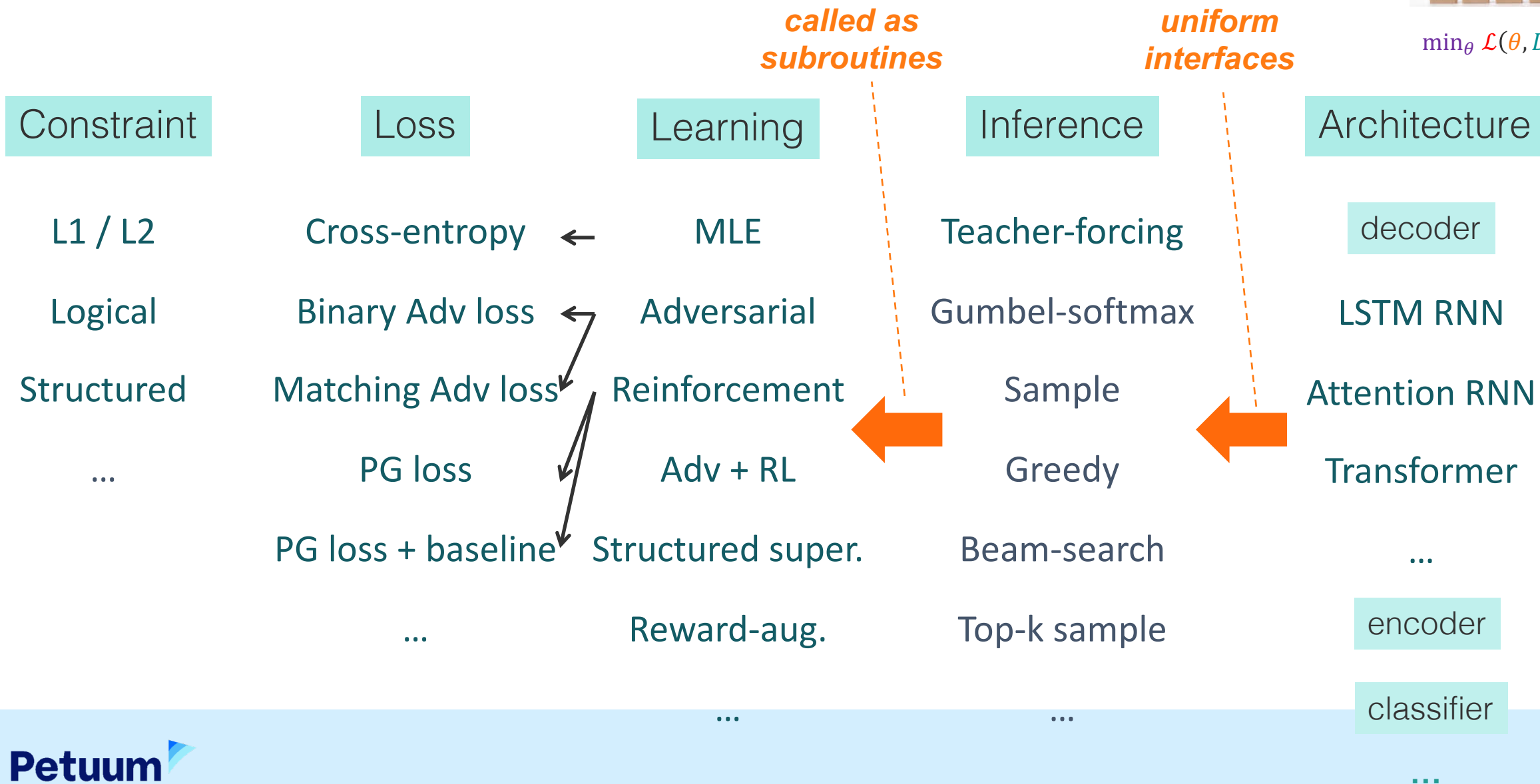
$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$



Holistic View

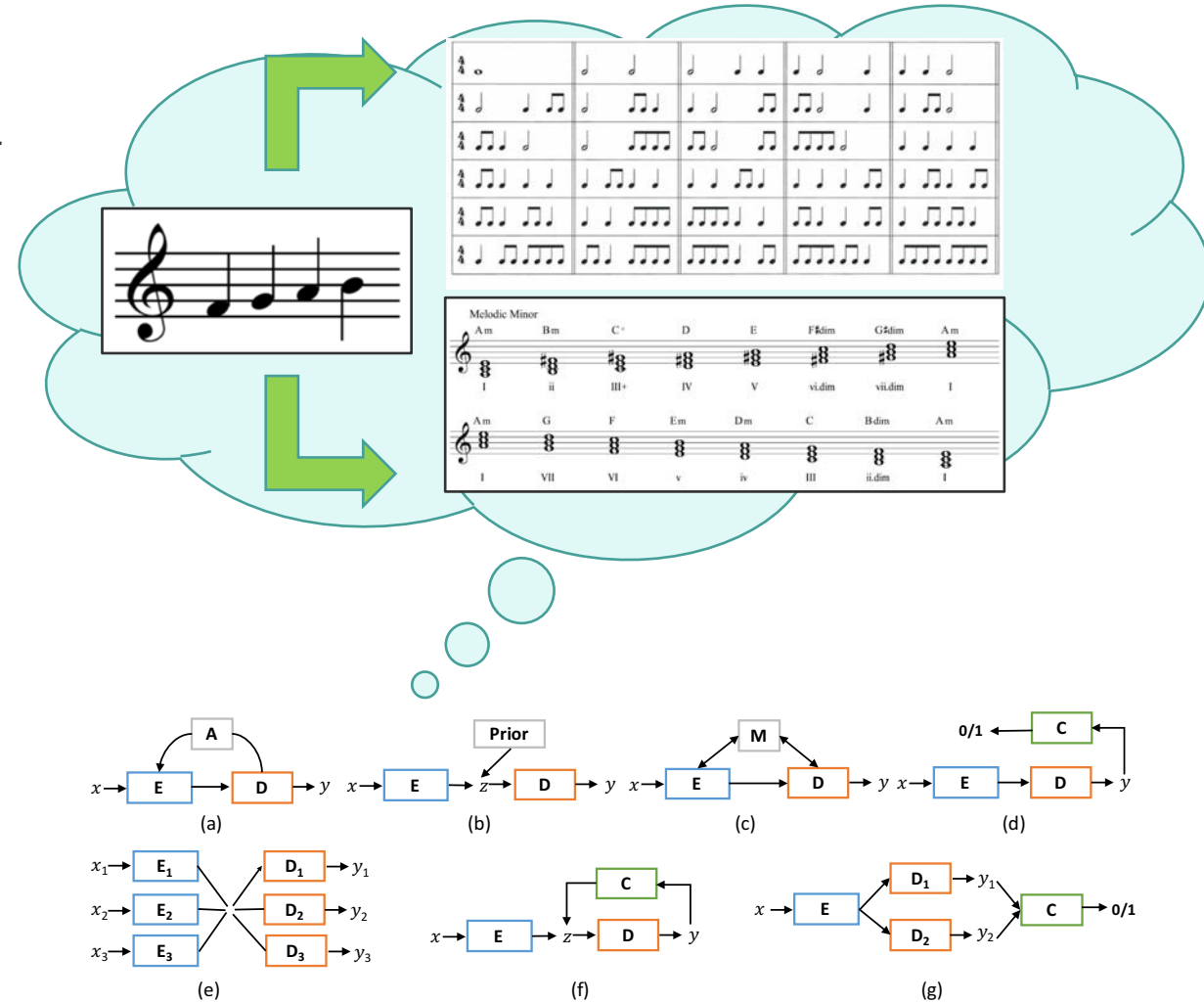


$$\min_{\theta} \mathcal{L}(\theta, D) + \Omega(\theta)$$

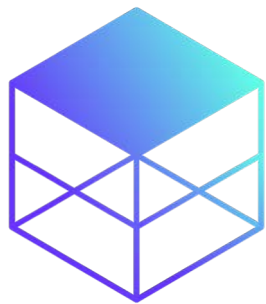


Composable ML – Take-Home Message

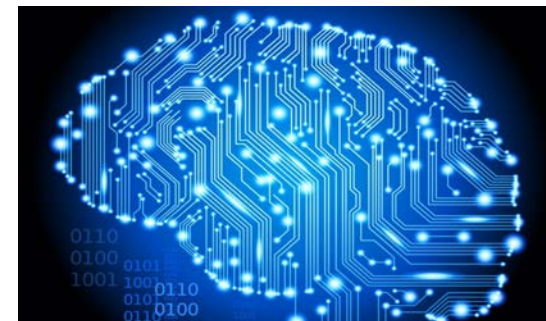
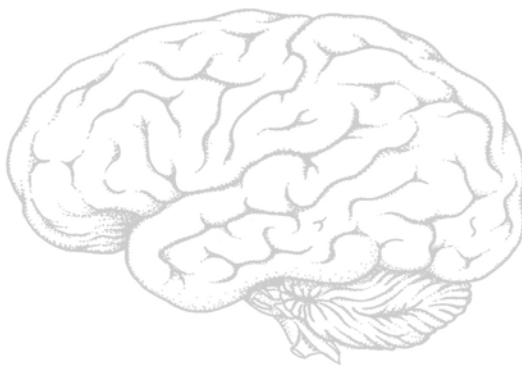
- Composable ML
 - Basic “musical notes” for complex ML systems
 - Structure/rules for combining notes into “chords”/“rhythms” ...
 - ... and chords/rhythms into compositions
 - (or just think of it as Lego for ML)
- Today – Symbolic programming via Petuum Texar open source
- Future – Drag and Play graphical UI



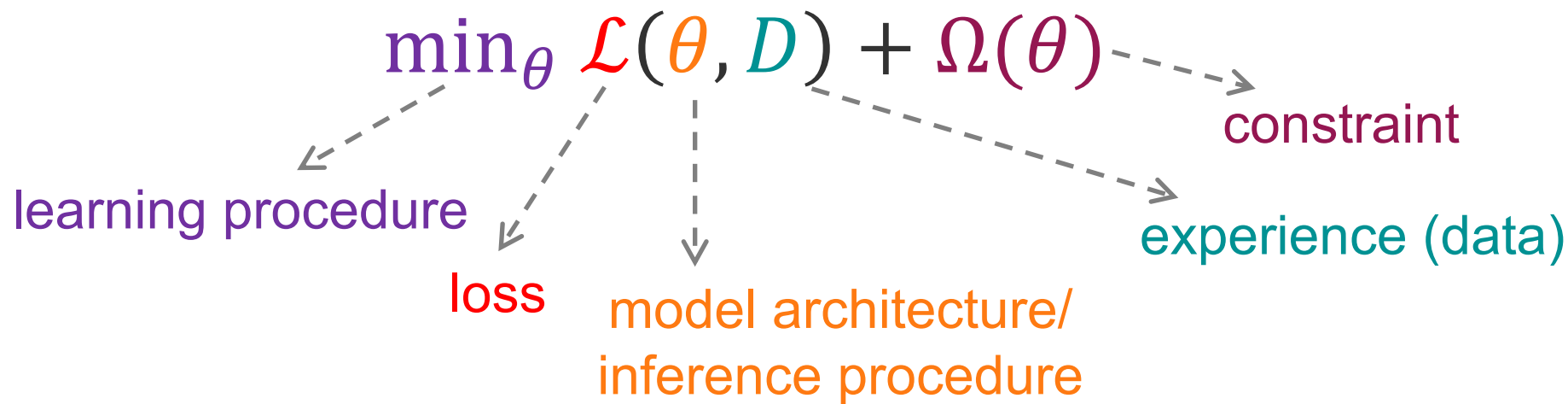
Texar



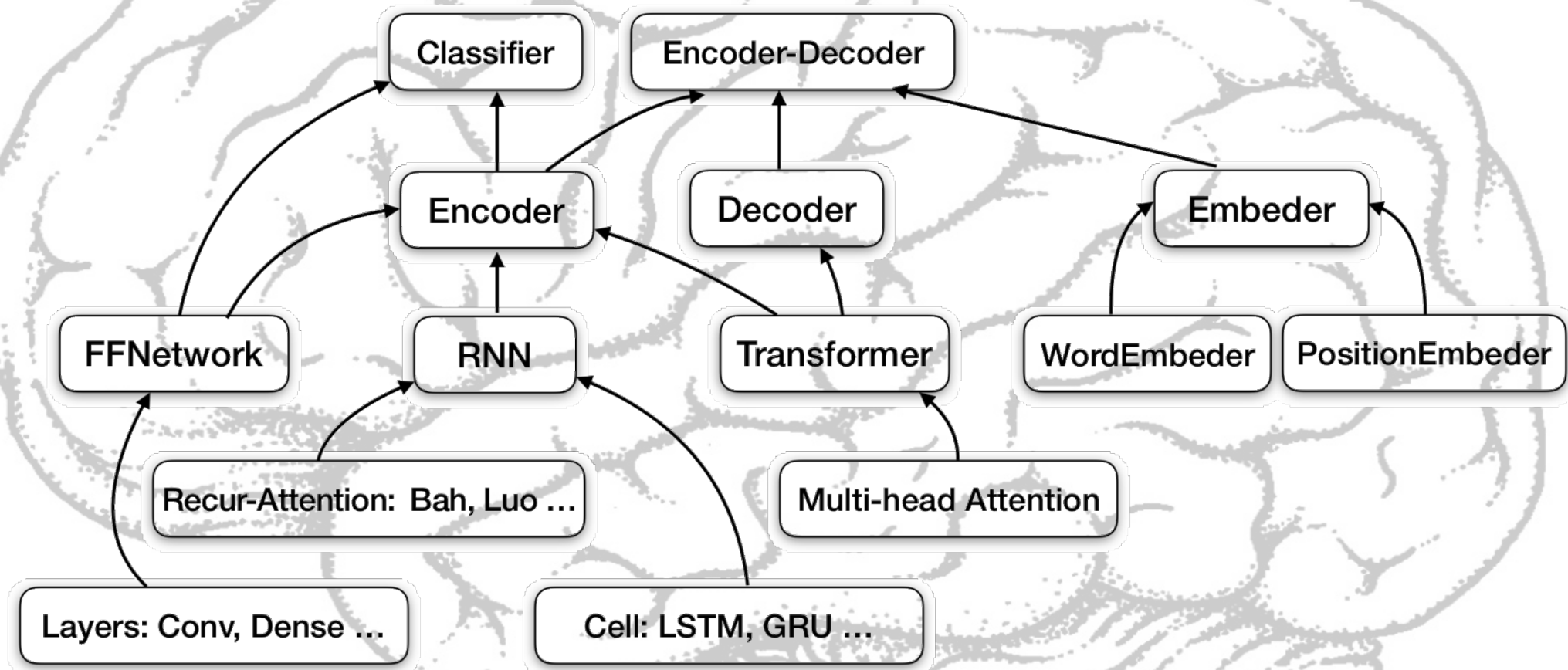
Compose your ML applications
like playing building blocks



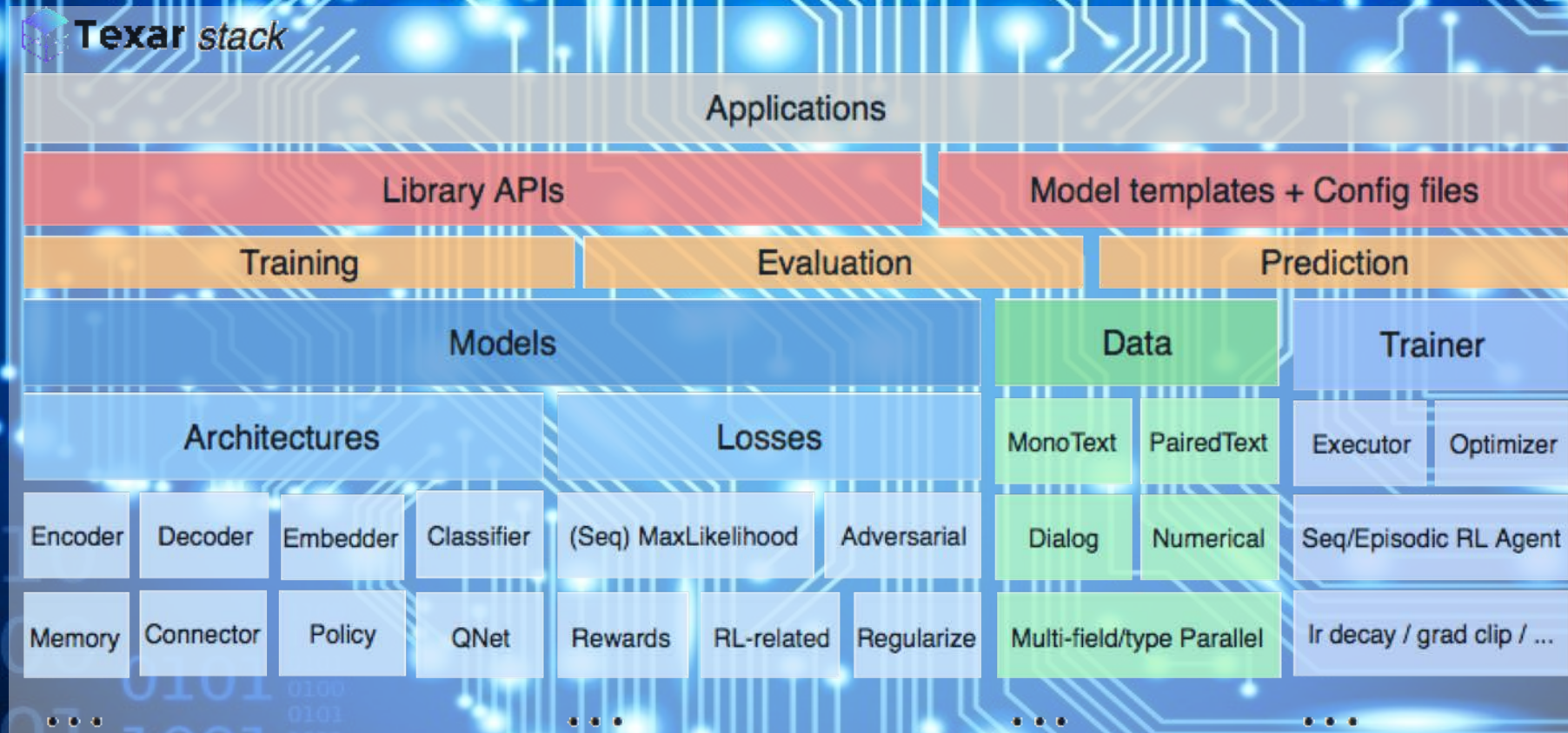
Decomposing Machine Learning



Expert's Intellectual "View" of Composable ML

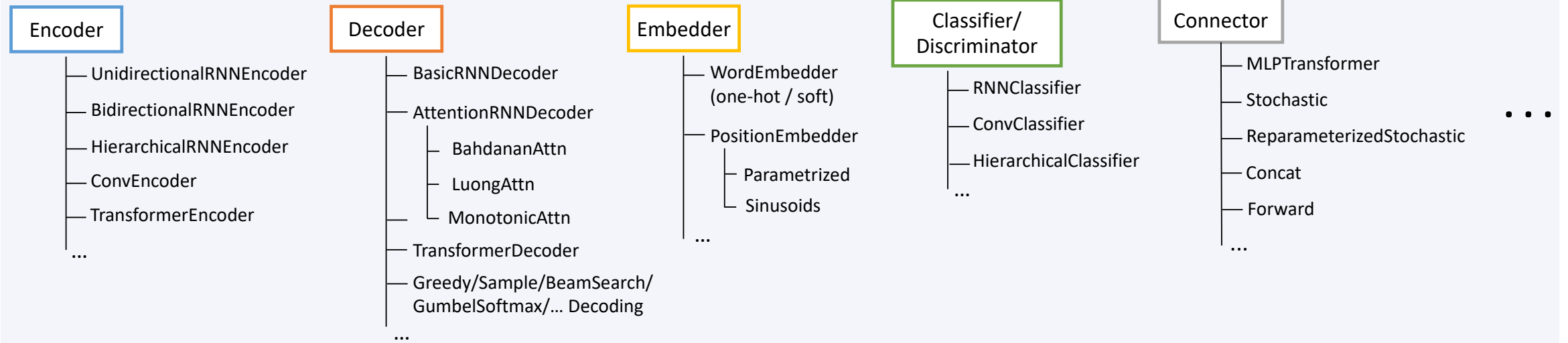


Texar Stack – Operationalized “View” of Composable ML

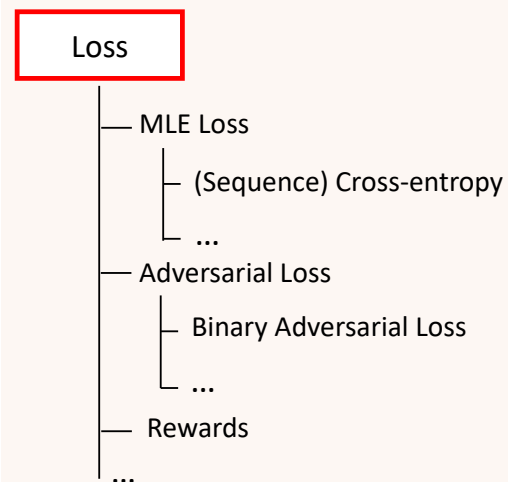


Module Catalog in Texar

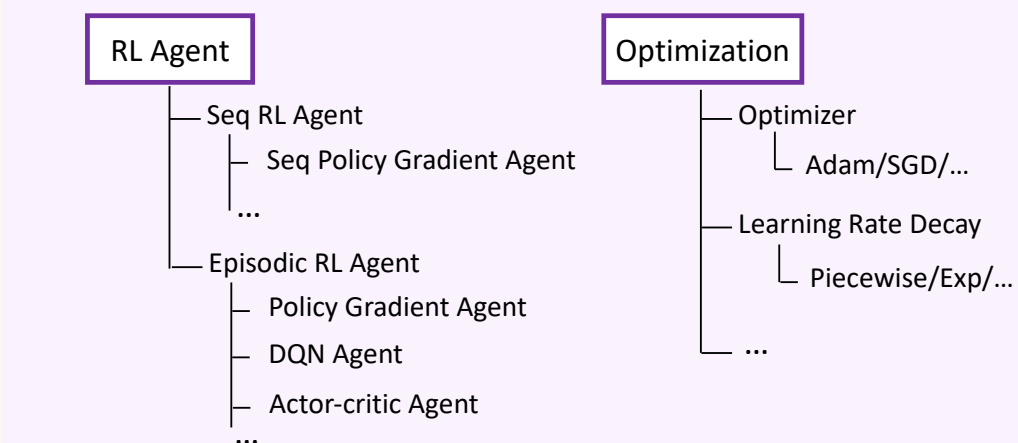
Model architecture



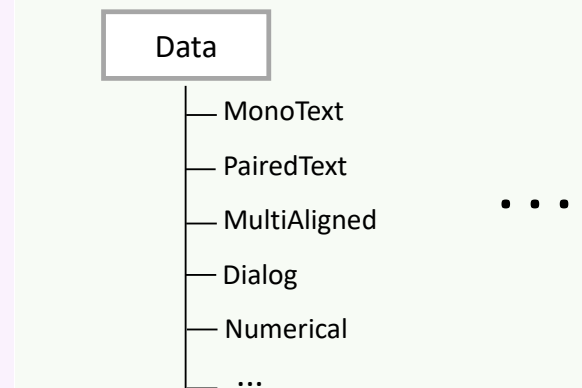
Model loss



Trainer



Data



Texar Highlights



Modularized

Assembles any complex model like playing building blocks



Versatile

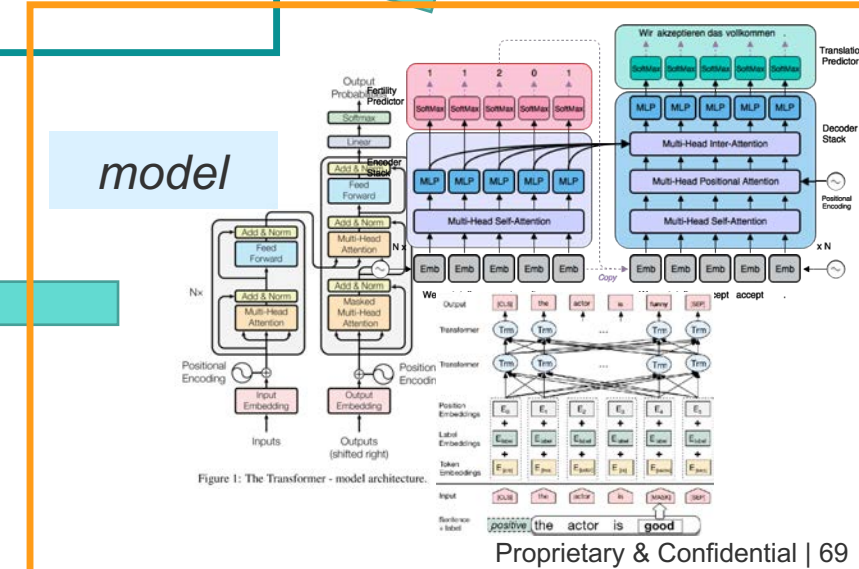
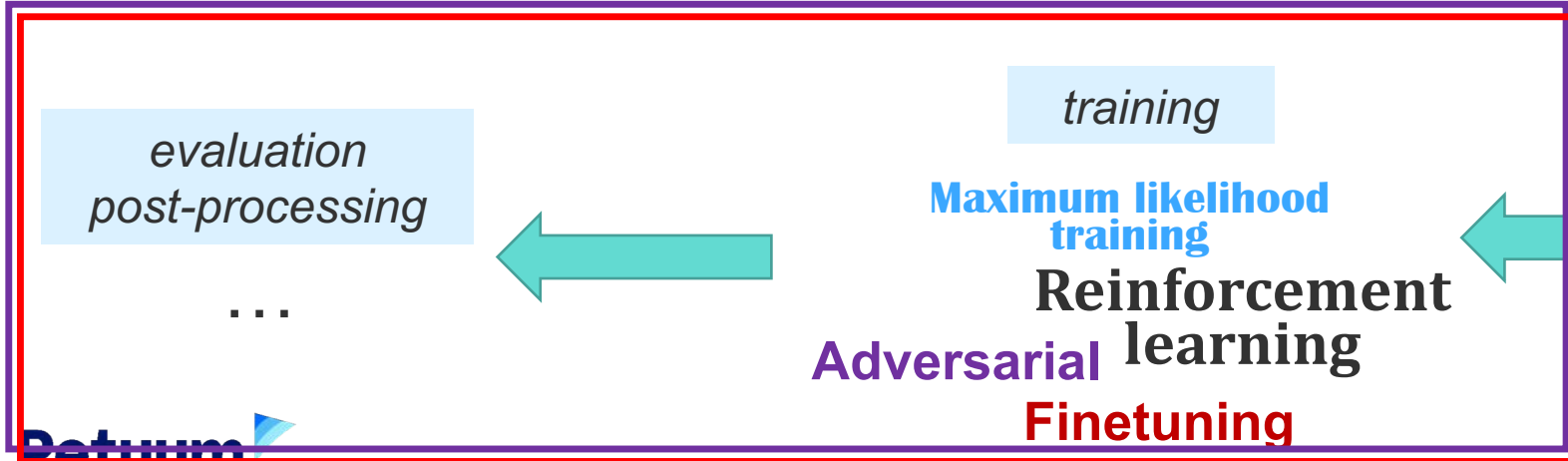
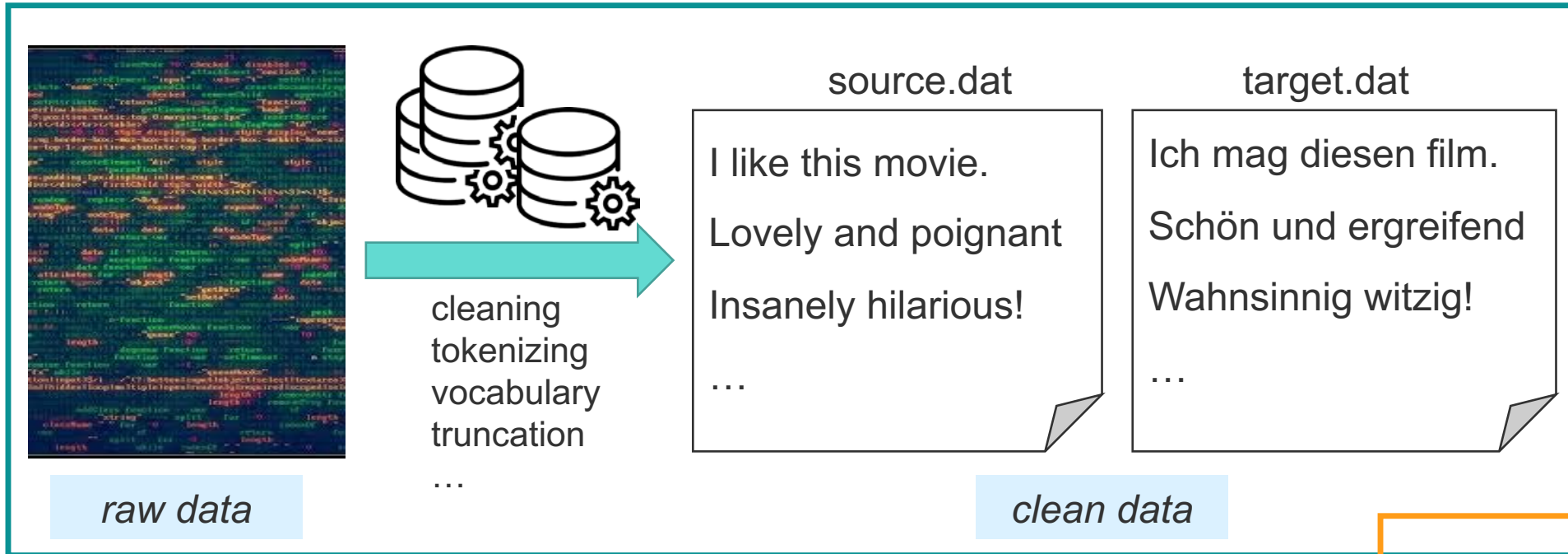
Supports a large variety of models/algorithms/applications ...



Extensible

Allows to plug in any customized or external modules

Running Example: Machine Translation



Running Example: Machine Translation - Data Preparation

- Input: a source sentence:

I like this movie.

- Output: a target sentence:

Ich mag diesen film.

- Dataset:

source.txt

I like this movie.
Lovely and poignant
Insanely hilarious!
...

target.txt

Ich mag diesen film.
Schön und ergreifend
Wahnsinnig witzig!
...

vocab.txt

I
like
this
movie
Ich
mag

Running Example: Machine Translation - Programming

Running Example: Machine Translation - Programming

Data

```
1 # Read data
2 dataset = PairedTextData(data_hparams)
3 batch = Dataloader(dataset).get_next()
```

```
{
  'batch_size': 64,
  'num_epochs': 10,
  'shuffle': True,
  'source_dataset': {
    'files': 'source.txt',
    'vocab_file': 'vocabulary.txt'
    'max_seq_length': 100,
    'bos_token': '<BOS>',
    'eos_token': '<EOS>',
    'embedding_init': { ... }
  },
  'target_dataset': {
    'files': 'target.txt',
  },
  # ...
}
```


Running Example: Machine Translation - Programming

Data

```
1 # Read data
2 dataset = PairedTextData(data_hparams)
3 batch = Dataloader(dataset).get_next()
```

Architecture
& Inference

Running Example: Machine Translation - Programming

Data

```
1 # Read data
2 dataset = PairedTextData(data_hparams)
3 batch = DataIterator(dataset).get_next()
4 # Encode
5 embedder = WordEmbedder(dataset.vocab.size, hparams=embedder_hparams)
```

Architecture
& Inference

```
{
  embedding_dim
  dropout
  initialization
  ...
}
```

Running Example: Machine Translation - Programming

Data

```
1 # Read data
2 dataset = PairedTextData(data_hparams)
3 batch = DataIterator(dataset).get_next()
4 # Encode
5 embedder = WordEmbedder(dataset.vocab.size, hparams=embedder_hparams)
6 encoder = TransformerEncoder(hparams=encoder_hparams)
```

Architecture
& Inference

```
{
  #blocks
  #heads
  hidden dim
  output dim
  dropout
  initialization
}
```

Running Example: Machine Translation - Programming

Data

```
1 # Read data
2 dataset = PairedTextData(data_hparams)
3 batch = Dataliterator(dataset).get_next()
4 # Encode
5 embedder = WordEmbedder(dataset.vocab.size, hparams=embedder_hparams)
6 encoder = TransformerEncoder(hparams=encoder_hparams)
7 enc_outputs = encoder(embedder(batch['source_text_ids']),
8                       batch['source_length'])
```

Architecture
& Inference

Running Example: Machine Translation - Programming

Data

```
1 # Read data
2 dataset = PairedTextData(data_hparams)
3 batch = Dataloader(dataset).get_next()
4 # Encode
5 embedder = WordEmbedder(dataset.vocab.size, hparams=embedder_hparams)
6 encoder = TransformerEncoder(hparams=encoder_hparams)
7 enc_outputs = encoder(embedder(batch['source_text_ids']),
8                       batch['source_length'])
9 # Build decoder
10 decoder = AttentionRNNDecoder(memory=enc_outputs,
11                               hparams=decoder_hparams)
12 # Maximum Likelihood Estimation
13 ## Teacher-forcing decoding
14 outputs, length, _ = decoder(decoding_strategy='teacher-forcing',
15                              inputs=embedder(batch['target_text_ids']),
16                              seq_length=batch['target_length']-1)
```

Architecture
& Inference

Running Example: Machine Translation - Programming

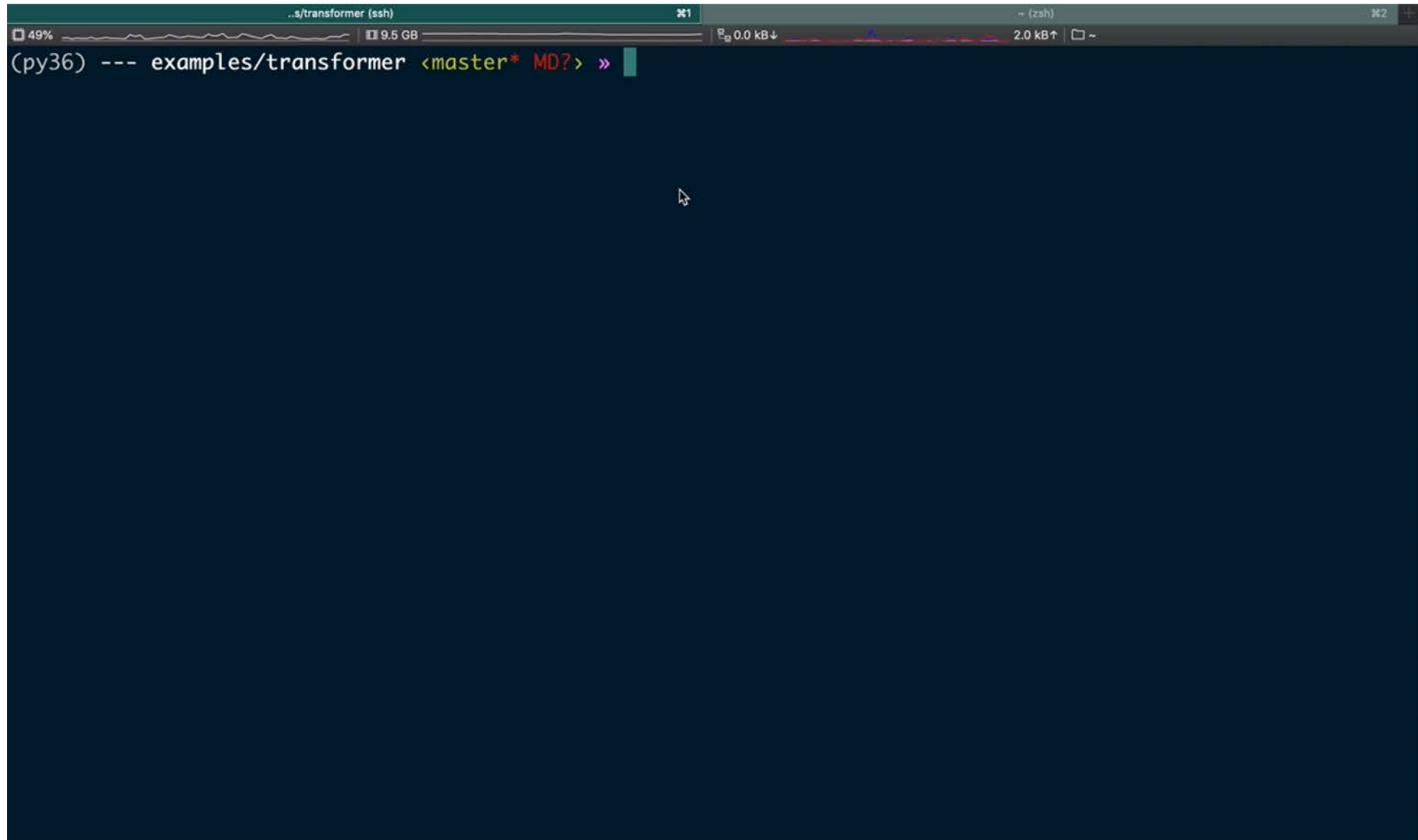
```
1 # Read data
2 dataset = PairedTextData(data_hparams)
3 batch = Datalerator(dataset).get_next()
4 # Encode
5 embedder = WordEmbedder(dataset.vocab.size, hparams=embedder_hparams)
6 encoder = TransformerEncoder(hparams=encoder_hparams)
7 enc_outputs = encoder(embedder(batch['source_text_ids']),
8                       batch['source_length'])
9 # Build decoder
10 decoder = AttentionRNNDecoder(memory=enc_outputs,
11                               hparams=decoder_hparams)
12 # Maximum Likelihood Estimation
13 ## Teacher-forcing decoding
14 outputs, length, _ = decoder(decoding_strategy='teacher-forcing',
15                             inputs=embedder(batch['target_text_ids']),
16                             seq_length=batch['target_length']-1)
17 ## Cross-entropy loss
18 loss = sequence_sparse_softmax_cross_entropy(
19     labels=batch['target_text_ids'][:,1:], logits=outputs.logits, seq_length=length)
20
```

Data

Architecture & Inference

Learning Loss

DEMO



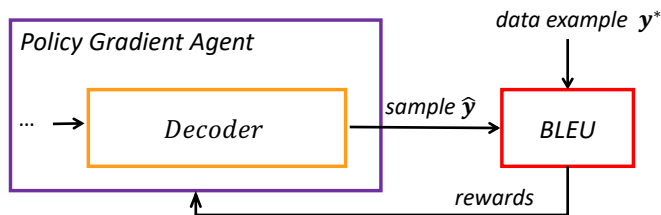
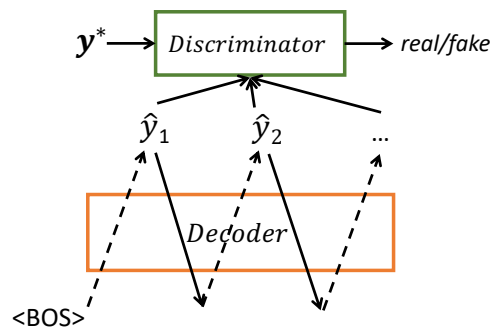
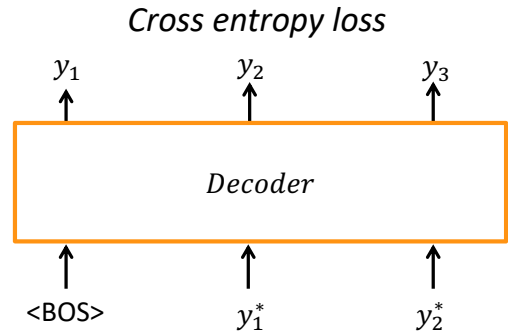
The image shows a terminal window with a dark blue background. At the top, there are two tabs: the first is titled "...s/transformer (ssh)" and the second is titled "~ (zsh)". Below the tabs, there is a status bar with various indicators: a battery icon at 49%, a memory usage indicator at 9.5 GB, and network activity indicators showing 0.0 kB down and 2.0 kB up. The main terminal area contains the prompt "(py36) --- examples/transformer <master* MD?> »" followed by a teal cursor. A mouse cursor is visible in the center of the terminal area.

Example (cont'd): Machine Translation

```
1 # Read data
2 dataset = PairedTextData(data_hparams)
3 batch = Dataloader(dataset).get_next()
4 # Encode
5 embedder = WordEmbedder(dataset.vocab.size, hparams=embedder_hparams)
6 encoder = TransformerEncoder(hparams=encoder_hparams)
7 enc_outputs = encoder(embedder(batch['source_text_ids']),
8                       batch['source_length'])
9 # Build decoder
10 decoder = AttentionRNNDecoder(memory=enc_outputs,
11                               hparams=decoder_hparams)
12 # Maximum Likelihood Estimation
13 ## Teacher-forcing decoding
14 outputs, length, _ = decoder(decoding_strategy='teacher-forcing',
15                             inputs=embedder(batch['target_text_ids']),
16                             seq_length=batch['target_length']-1)
17 ## Cross-entropy loss
18 loss = sequence_sparse_softmax_cross_entropy(
19     labels=batch['target_text_ids'][:,1:], logits=outputs.logits, seq_length=length)
20
```

Maximum likelihood Estimation

Different Learning Algorithms



- Maximum Likelihood Estimation
- Adversarial Learning
- Reinforcement Learning

Switching between Learning Algorithms

Data

```
1 # Read data
2 dataset = PairedTextData(data_hparams)
3 batch = Datalerator(dataset).get_next()
```

Keep unchanged

```
4 # Encode
5 embedder = WordEmbedder(dataset.vocab.size, hparams=embedder_hparams)
6 encoder = TransformerEncoder(hparams=encoder_hparams)
7 enc_outputs = encoder(embedder(batch['source_text_ids']),
8                       batch['source_length'])
```

Architecture
& Inference

```
9 # Build decoder
10 decoder = AttentionRNNDecoder(memory=enc_outputs,
11                               hparams=decoder_hparams)
```

```
12 # Maximum Likelihood Estimation
13 ## Teacher-forcing decoding
14 outputs, length, _ = decoder(decoding_strategy='teacher-forcing',
15                             inputs=embedder(batch['target_text_ids']),
16                             seq_length=batch['target_length']-1)
```

**Maximum likelihood
Estimation**

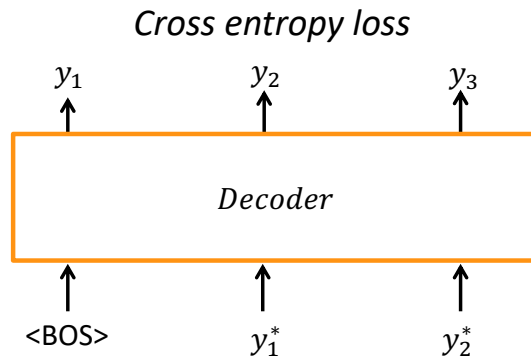
Learning
Loss

```
17 ## Cross-entropy loss
18 loss = sequence_sparse_softmax_cross_entropy(
19     labels=batch['target_text_ids'][:,1:], logits=outputs.logits, seq_length=length)
```

20

Switching from MLE to Adversarial Learning

- Maximum likelihood



```
# Teacher-forcing decoding
```

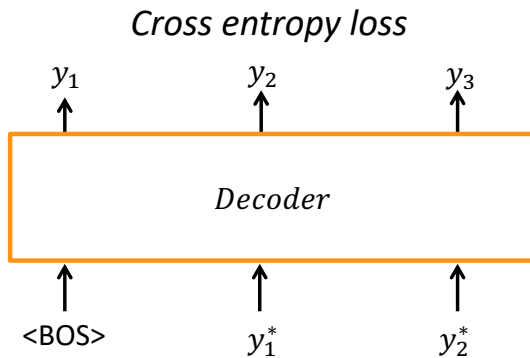
```
outputs, length, _ = decoder(decoding_strategy='teacher-forcing',  
                             inputs=embedder(batch['target_text_ids']),  
                             seq_length=batch['target_length']-1)
```

```
# Cross-entropy loss
```

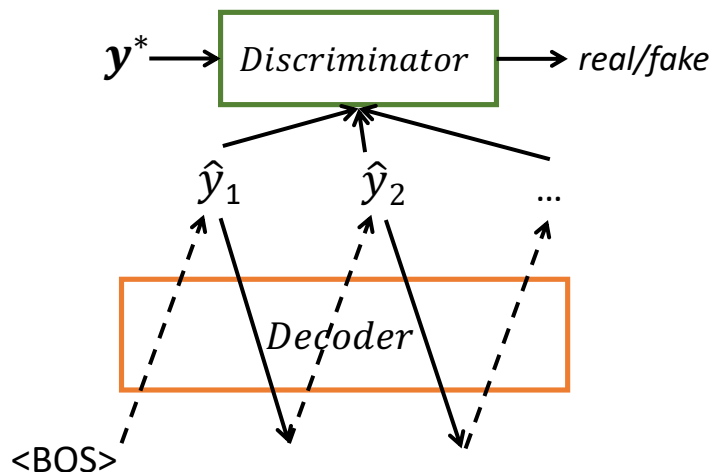
```
loss = sequence_sparse_softmax_cross_entropy(  
    labels=batch['target_text_ids'][:,1:], logits=outputs.logits, seq_length=length)
```

Switching from MLE to Adversarial Learning

- Maximum likelihood



- Adversarial learning



```
# Teacher-forcing decoding
```

```
outputs, length, _ = decoder(decoding_strategy='teacher-forcing',  
                             inputs=embedder(batch['target_text_ids']),  
                             seq_length=batch['target_length']-1)
```

```
# Cross-entropy loss
```

```
loss = sequence_sparse_softmax_cross_entropy(  
    labels=batch['target_text_ids'][:,1:], logits=outputs.logits, seq_length=length)
```

```
# Gumbel-softmax decoding
```

```
helper = GumbelSoftmaxTrainingHelper(  
    start_tokens=[BOS]*batch_size, end_token=EOS, embedding=embedder)  
outputs, _, _ = decoder(helper=helper)
```

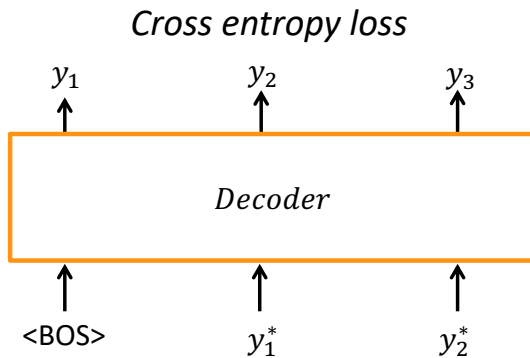
```
discriminator = Conv1DClassifier(hparams=conv_hparams)
```

```
# Binary adversarial loss
```

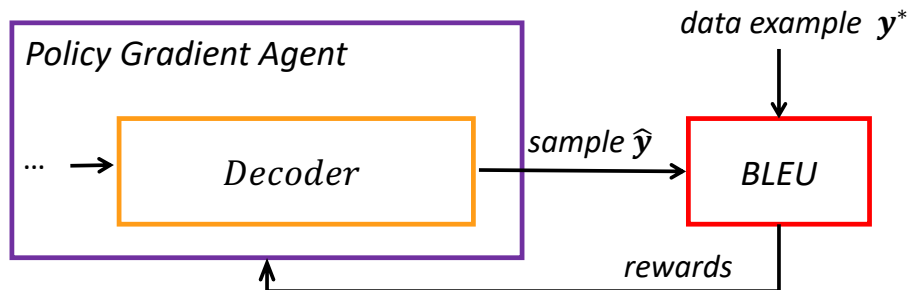
```
G_loss, D_loss = binary_adversarial_losses(  
    embedder(batch['target_text_ids'][:, 1:]),  
    embedder(soft_ids=softmax(outputs.logits)),  
    discriminator)
```

Switching from MLE to Reinforcement Learning

- Maximum likelihood



- Reinforcement learning



```
# Teacher-forcing decoding
```

```
outputs, length, _ = decoder(decoding_strategy='teacher-forcing',
                              inputs=embedder(batch['target_text_ids']),
                              seq_length=batch['target_length']-1)
```

```
# Cross-entropy loss
```

```
loss = sequence_sparse_softmax_cross_entropy(
    labels=batch['target_text_ids'][:,1:], logits=outputs.logits, seq_length=length)
```

```
# Random sample decoding
```

```
outputs, length, _ = decoder(decoding_strategy='random_sample',
                              start_tokens=[BOS]*batch_size, end_token=EOS,
                              embedding=embedder)
```

```
# Policy gradient agent for learning
```

```
agent = SeqPGAgent(
    samples=outputs.sample_id, logits=outputs.logits, seq_length=length)
```

```
for _ in range(STEPS):
```

```
    samples = agent.get_samples()
```

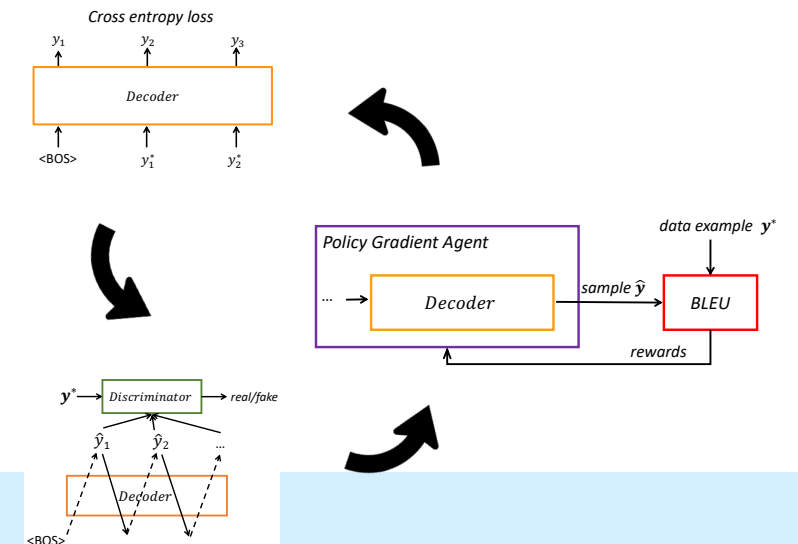
```
    rewards = BLEU(batch['target_text_ids'], samples) # Reward
```

```
    agent.observe(rewards)
```

Summary of MT in Texar

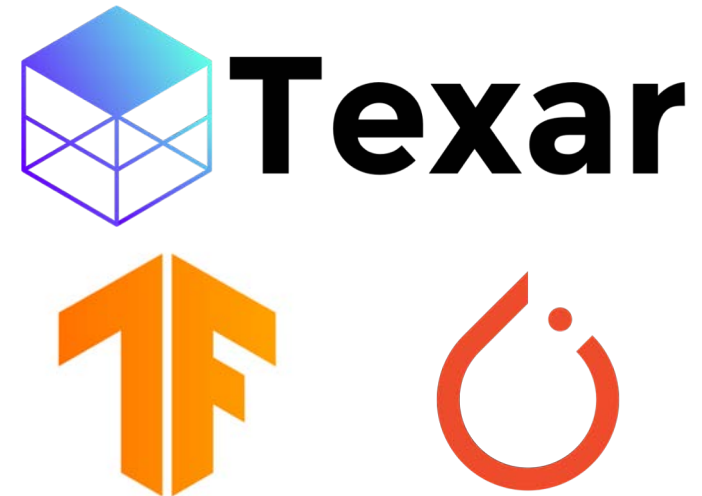
- Highly modularized programming
 - Data, architecture, loss, inference, learning, ...
 - Intuitive conceptual-level APIs
- Easy switch between learning algorithms
 - Plug in & out modules
 - No changes to irrelevant parts

```
1 # Read data
2 dataset = PairedTextData(data_hparams)
3 batch = Dataliterator(dataset).get_next()
4 # Encode
5 embedder = WordEmbedder(dataset.vocab.size, hparams=embedder_hparams)
6 encoder = TransformerEncoder(hparams=encoder_hparams)
7 enc_outputs = encoder(embedder(batch['source_text_ids']),
8                       batch['source_length'])
9 # Build decoder
10 decoder = AttentionRNNDecoder(memory=enc_outputs,
11                               hparams=decoder_hparams)
12 # Maximum Likelihood Estimation
13 ## Teacher-forcing decoding
14 outputs, length, _ = decoder(decoding_strategy='teacher-forcing',
15                              inputs=embedder(batch['target_text_ids']),
16                              seq_length=batch['target_length']-1)
17 ## Cross-entropy loss
18 loss = sequence_sparse_softmax_cross_entropy(
19     labels=batch['target_text_ids'][:,1:], logits=outputs.logits, seq_length=length)
20
```

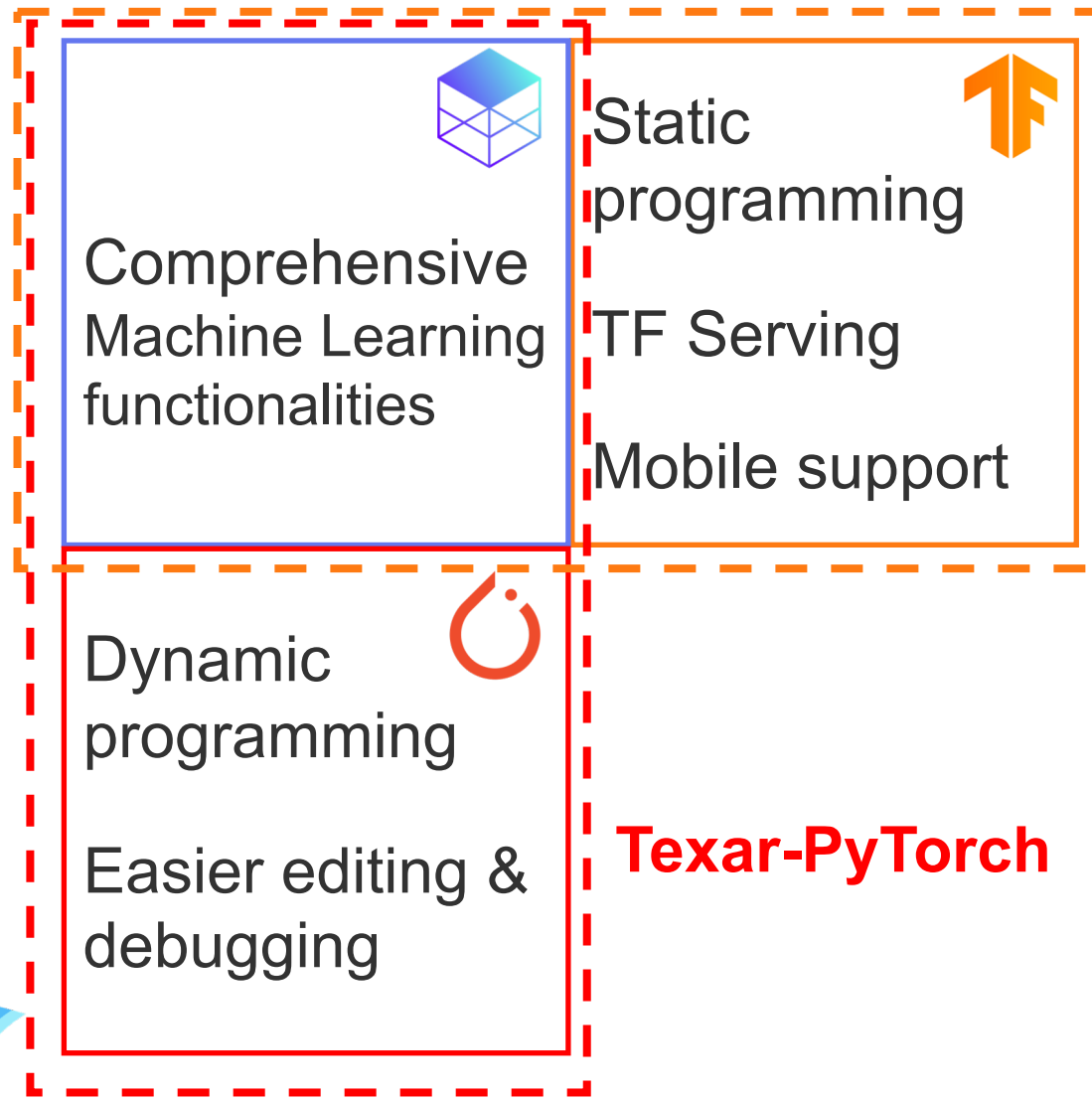


Other Features (I): Support of TensorFlow and PyTorch

- Texar is built upon TF and PyTorch
 - Texar-TF & Texar-PyTorch: mostly the same interfaces!
 - Higher-level intuitive APIs without loss of flexibility
 - Lots of ML components ready to use
- Combine the best design of TF and PyTorch
 - TF:
 - Easy and efficient data processing APIs
 - Excellent factorization of ML modules
 - Turnkey model training processor
 - PyTorch:
 - Intuitive programming interfaces
 - Transparent variable scope and sharing to users



Other Features (I): Support of TensorFlow and PyTorch



Texar-TF

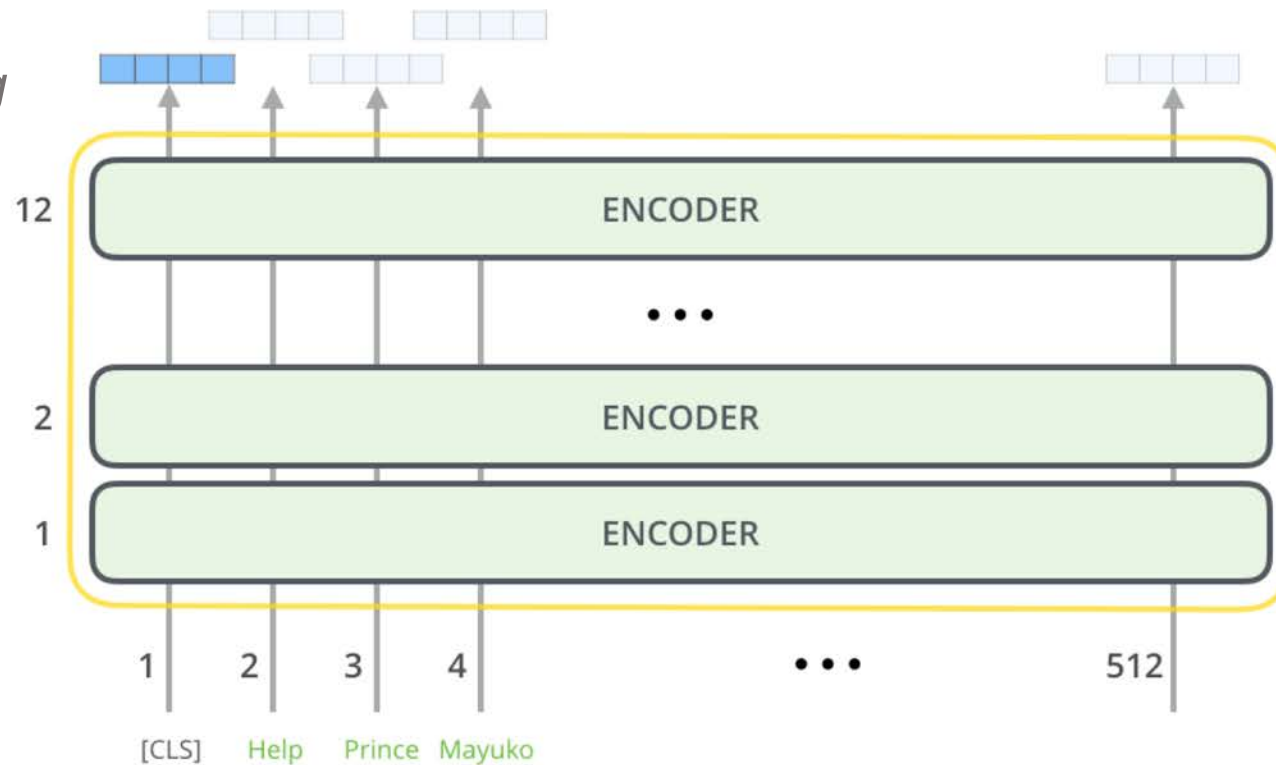
- Texar provides the same comprehensive supports of ML functionalities
- Fully compatible with native TF & PyTorch APIs
 - Get all other useful features of TF or PyTorch with **Texar-TF** & **Texar-PyTorch**

Other Features (II): SOTA Pretrained Models



Bert [Devlin et al., 2018]

contextual embedding



input sentence

Other Features (II): SOTA Pretrained Models



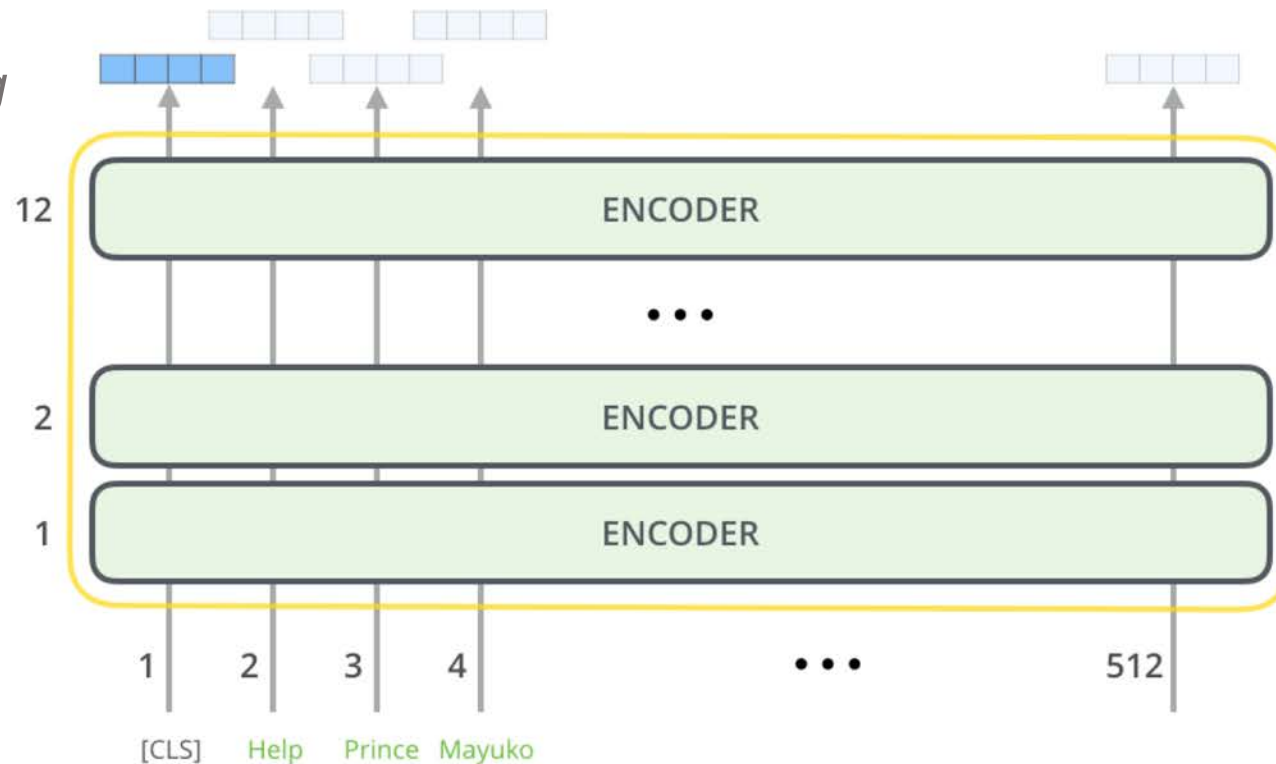
Bert [Devlin et al., 2018]

feature
extraction

```
model = BertEncoder()
```

```
features = model(input_ids, input_length, segment_ids)
```

contextual embedding



input sentence

Other Features (II): SOTA Pretrained Models



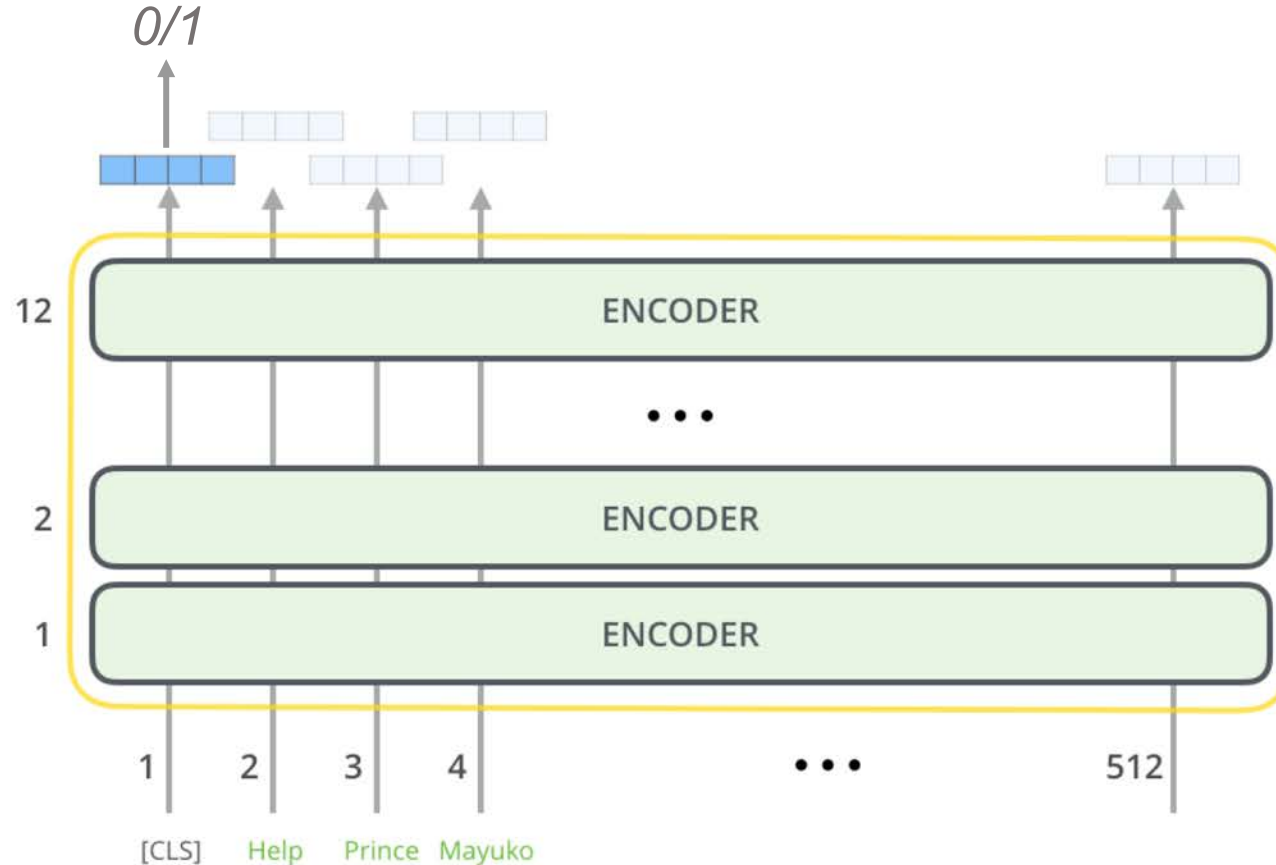
Bert [Devlin et al., 2018]

sequence
classification

```
model = BertClassifier(hparams={'clas_strategy': 'cls_time'})
```

```
logits, preds = model(input_ids, input_length, segment_ids)
```

sentence class



input sentence

Other Features (II): SOTA Pretrained Models



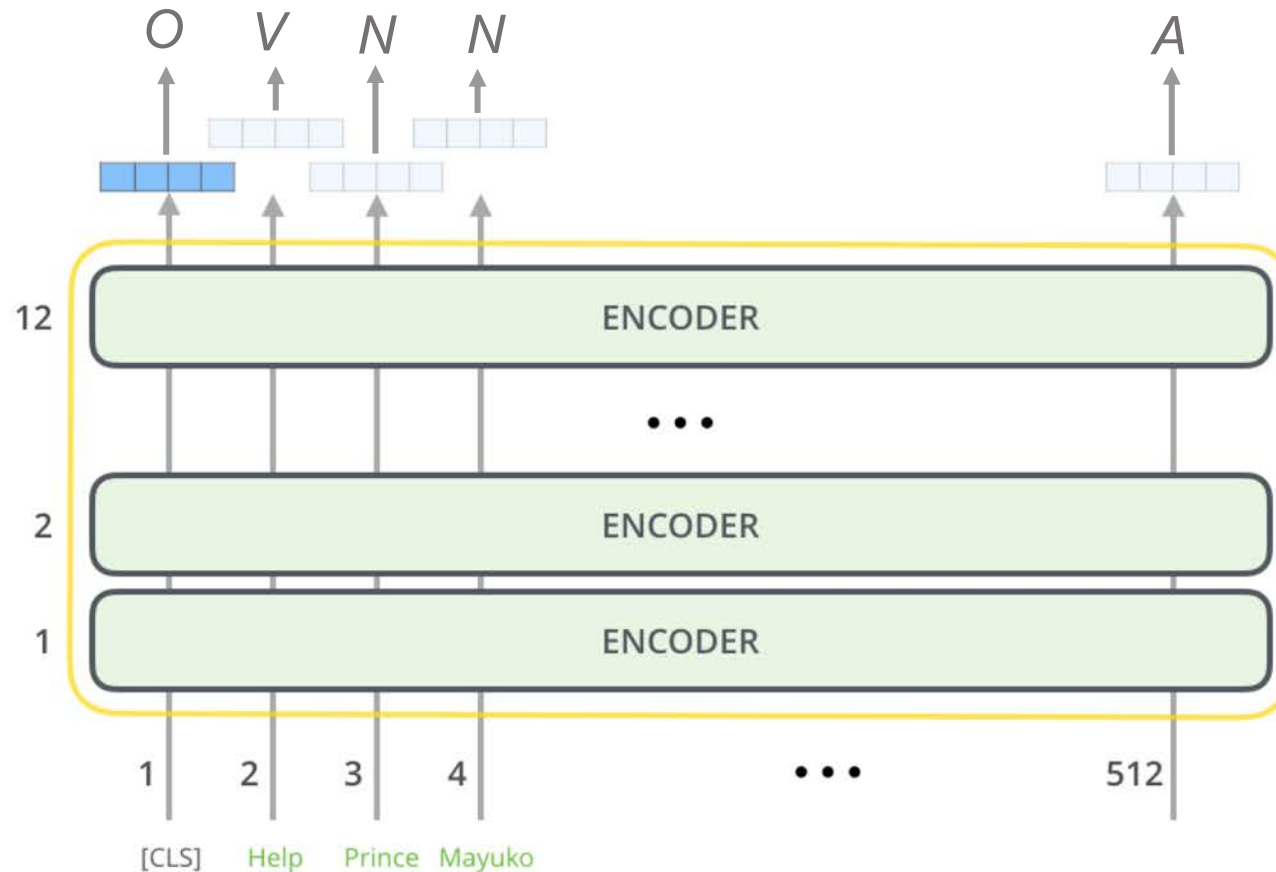
Bert [Devlin et al., 2018]

sequence
labeling

```
model = BertClassifier(hparams={'clas_strategy': 'all_time'})
```

```
logits, preds = model(input_ids, input_length, segment_ids)
```

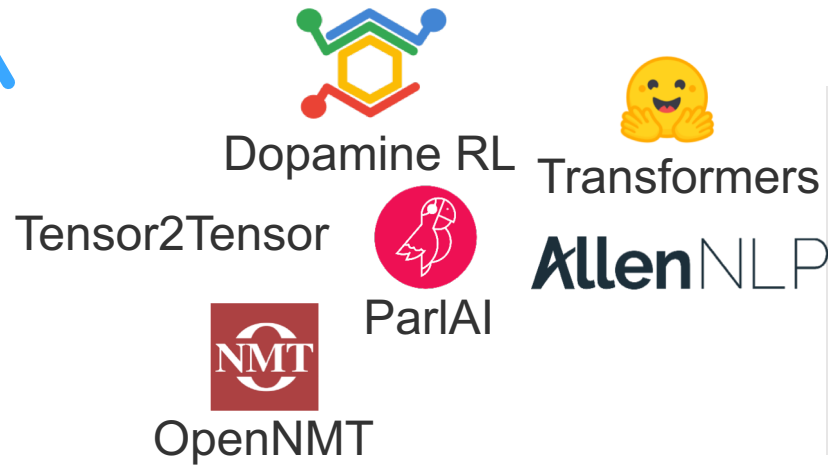
label per step



input sentence

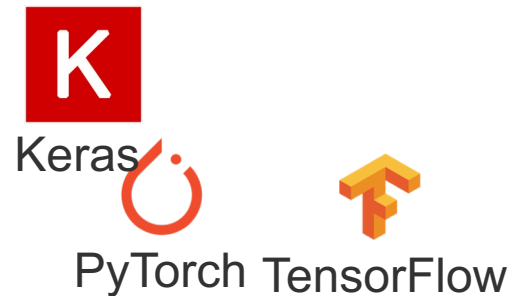
Spectrum of Existing Tools

*domain-specific,
higher-level APIs*



Interfaces at multiple abstraction levels

- Simplified APIs for common functionalities
- Advanced APIs for advanced functionalities and customizability



Spectrum of Existing Tools

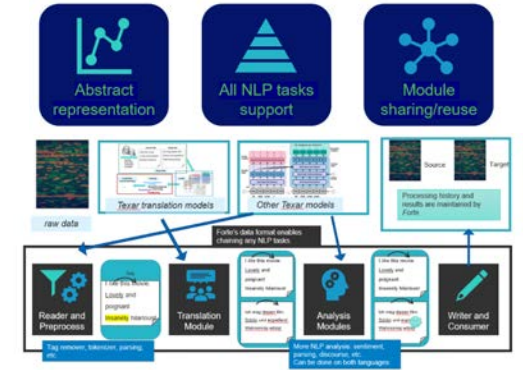
- Texar
 - Composable ML system – create ML models of complex design from ML first principles
- TensorFlow, PyTorch
 - Symbolic languages for low-level development of ML models
- BERT, GPT2, ...
 - ML models used to model human languages; also can refer to their downloadable pre-trained weights (ready-to-run)
- Amazon SageMaker, Google AutoML, Microsoft Azure ML
 - Commercial products for training, inference and management of pre-built ML models; may perform algorithm-driven model parameter or architecture tuning

Applications of Texar

Many products built on Texar

- ❑ FORTE – templates for larger complex NLP applications
- ❑ Chest X-Ray report writer
- ❑ Medical Registry report writer
- ❑ ICD coding system
- ❑ Financial knowledge base builder
- ❑ Financial summary/report writer
- ❑ Multi-Lingual Cognitive Chat Bots
 - ❑ For Call Center Support
 - ❑ For Retail In-Store Assistance

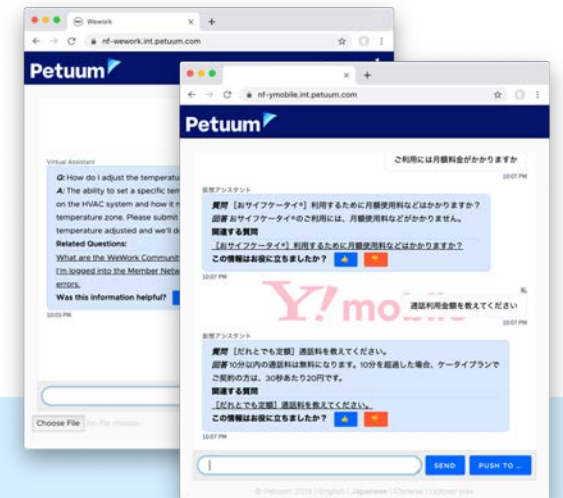
FORTE



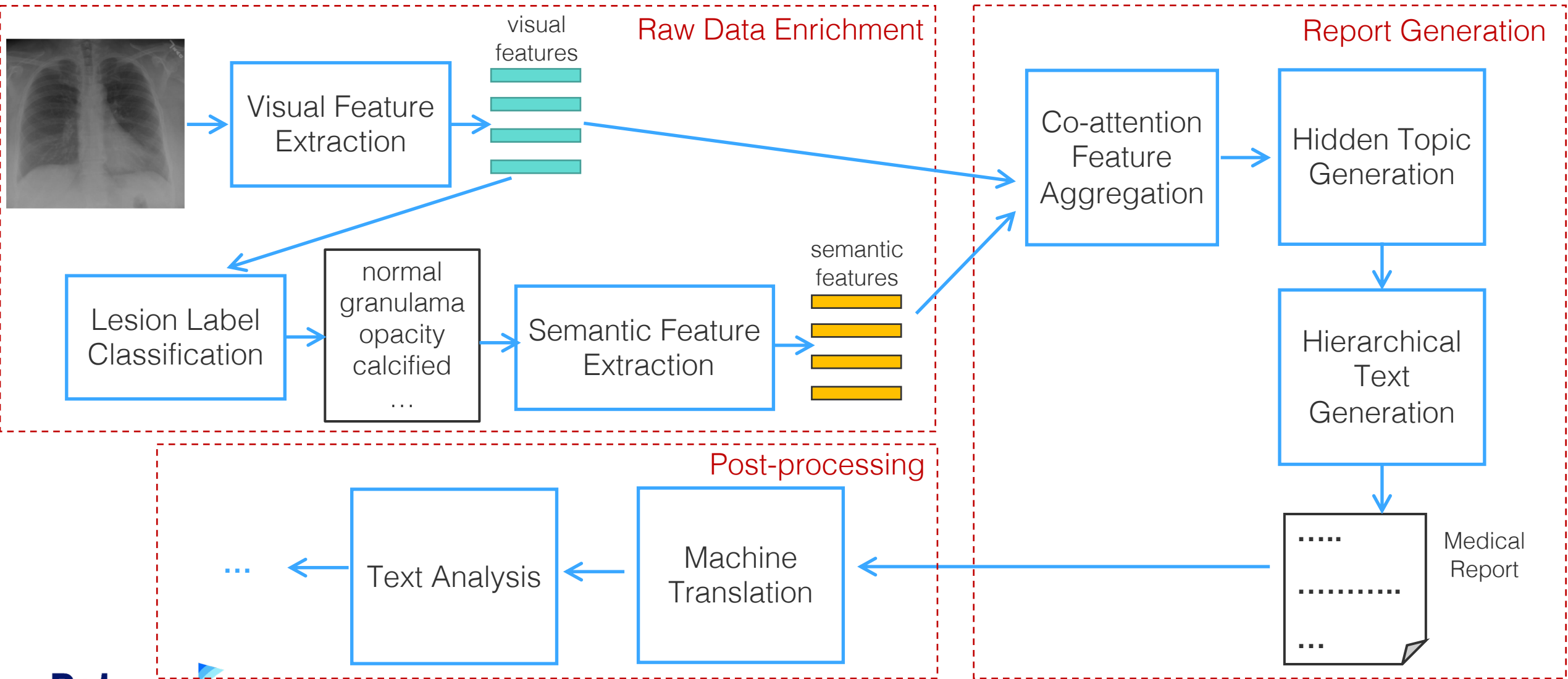
Chest X-Ray Report Writer



Multi-Lingual Cognitive Chat Bots



Example: Chest X-Ray Report Writer



Example: Chest X-Ray Report Writer + Translation

无胸部X光片
拖拽至此上传胸部X光片(.png)
上传照片

Texar Resources



Texar-TF

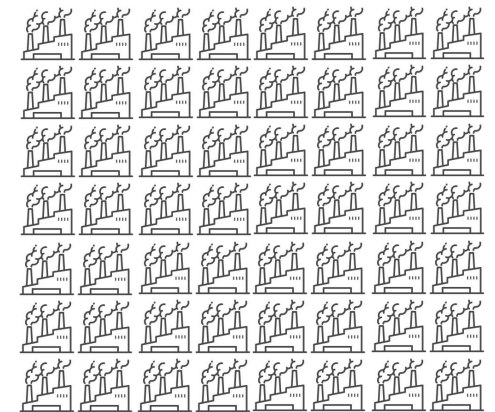


Texar-PyTorch

- Website: <https://asym1.io>
- GitHub (TF version): <https://github.com/asym1/texar>
- GitHub (PyTorch version): <https://github.com/asym1/texar-pytorch>
- Examples: <https://github.com/asym1/texar/blob/master/examples>
- Documentation: <https://texar.readthedocs.io/>
- Blog: <https://medium.com/@texar>
- Tech report: <https://arxiv.org/pdf/1809.00794.pdf>

Scalable AI Infrastructure

Scaling out AI applications –
easier, faster, cheaper



What is needed to scale AI?

Inter-Task Interfacing Application Templates

Assemble many tasks with uniform interfaces to form larger, complex applications

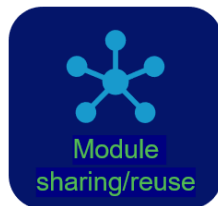
High Speed & Scalability

Made-to-order, just-in-time distributed training strategies

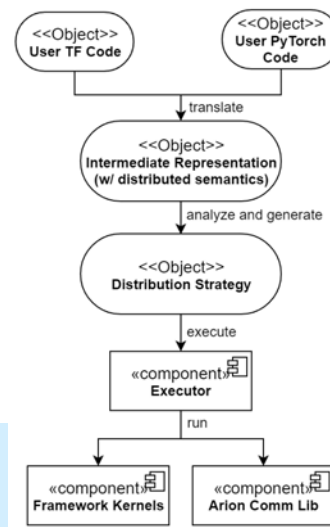
Save Cost & Time

Dynamically reassign CPU & GPU resources for fastest/cheapest workload completion time

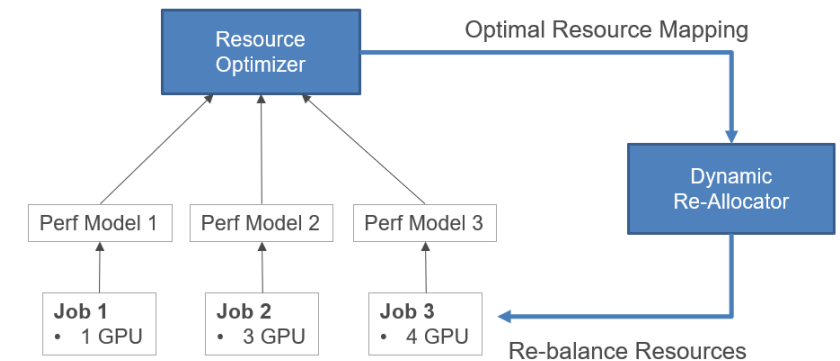
FORTE



ARION

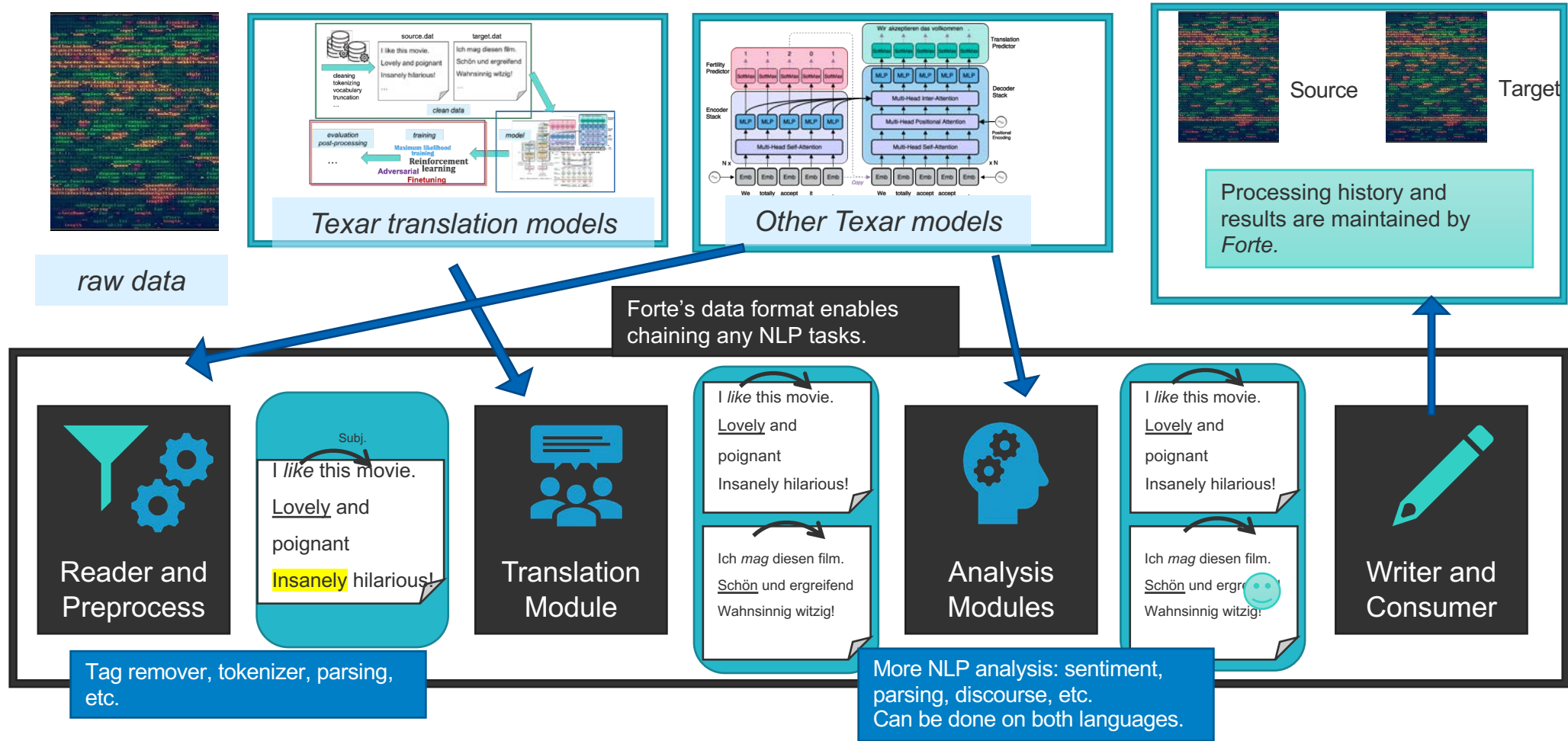


ESPER



Forte – Flexible Scaffold for Larger NLP Applications

Built upon *Texar*

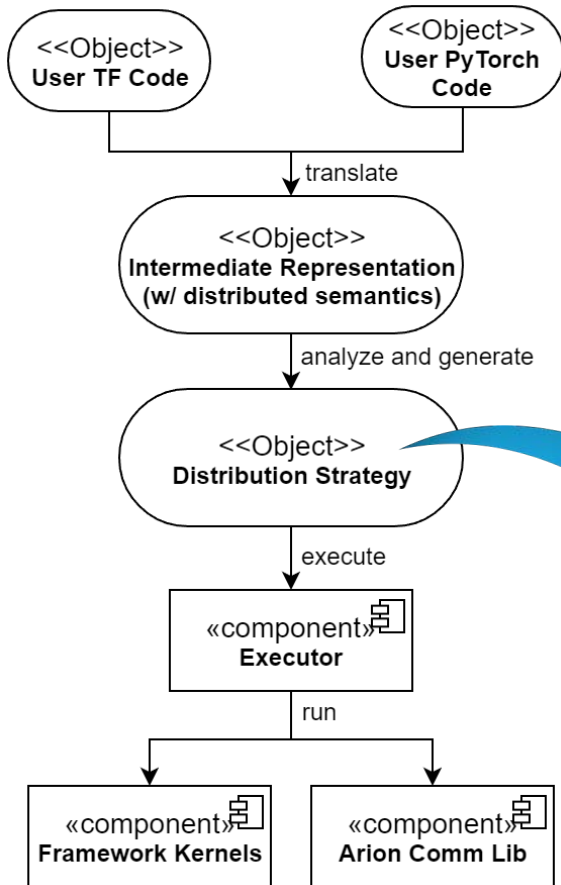


Arion – High Speed, Scalability for ML Model Training

Tailor-made distributed training strategies for hard-to-scale ML models

Pain Points

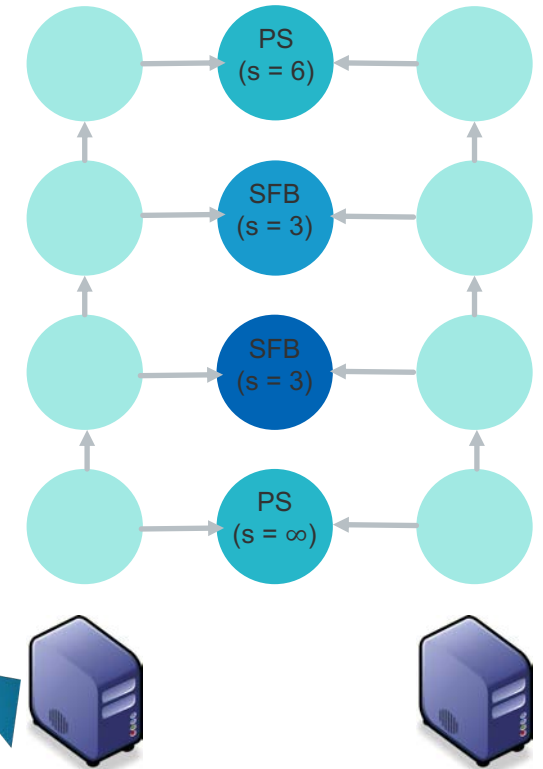
- Existing systems are mostly specialized and implemented based on 1 (probably at most 2) of those system architectures.
- | | | | | | |
|-----------------|--------------|-------------------|------------------|---------------|----------------------|
| AR | PS | SSP | SFB | Decentralized | Gradient Compression |
| Pipelining (NN) | Horovod (AR) | Poseidon (PS+SFB) | Dynamic Batching | | |
- Distribution performance/scalability is only achieved under certain condition, for a narrow family of models where the system is specialized.



Solution

Arion expresses models and training algorithms as an intermediate representation (IR) with distributed execution semantics

Arion rewrites IR and generates **tailor-made distributed execution strategy** using multiple system architectures (e.g. PS + SFB)

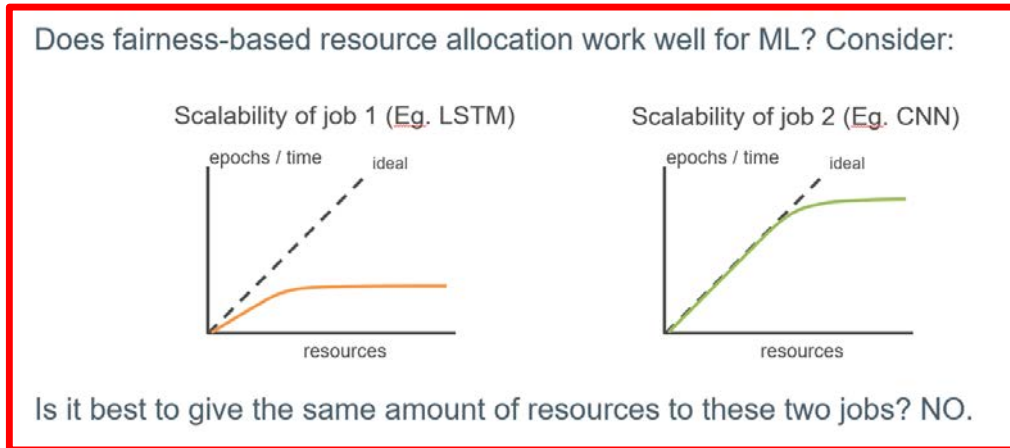
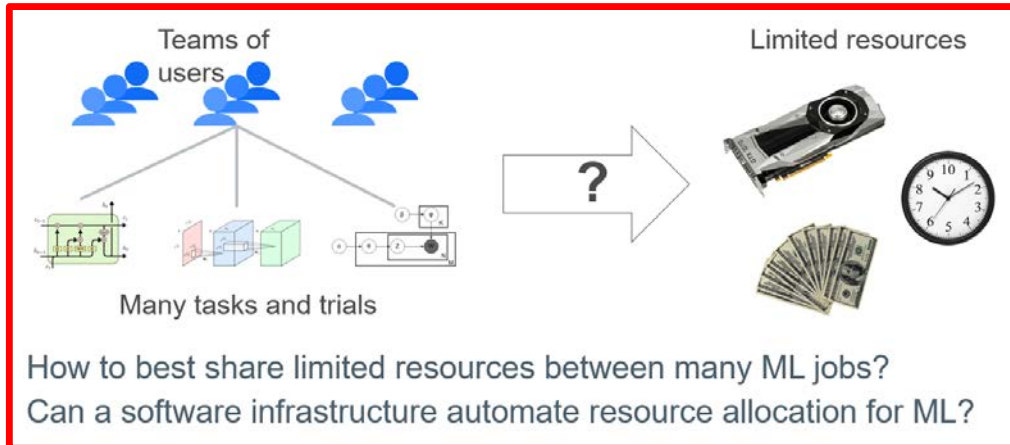


Value: Model training – from un-scalable (speed $\propto 1$)
to highly-scalable (speed $\propto n$ machines)

Esper – Save Compute Costs and Time

Dynamically reassign CPU/GPU to the most impactful work

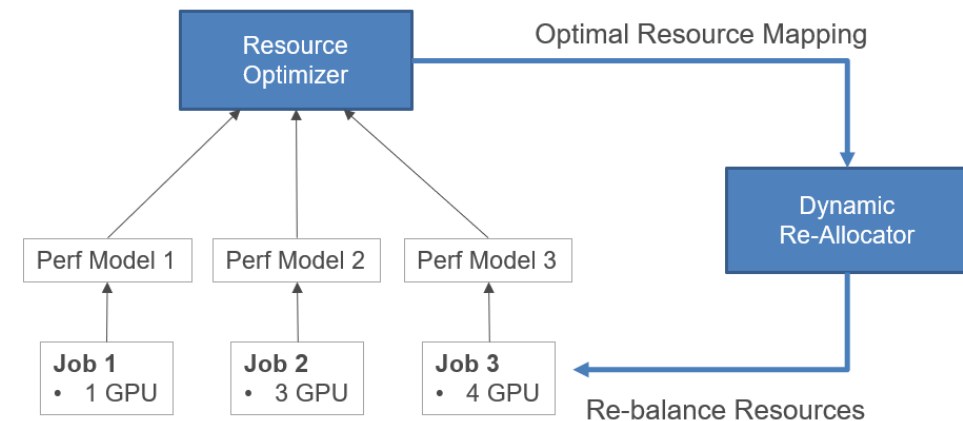
Pain Points



Solution

Esper learns the performance model for each model training job

Esper uses performance models to optimize (in near-real-time) CPUs/GPUs needed by each running job on the cluster/cloud



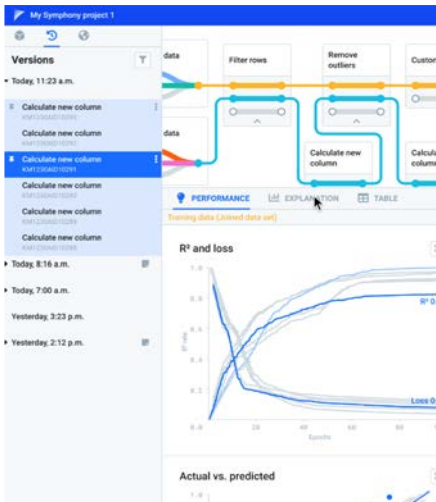
Value: Finish training *multiple* models faster (and cheaper!)

Esper is 2-4x faster vs Google Kubeflow

What is needed to productize AI?

Drag & Play UI
Manage, Experiment, Version
Systems, Tools, Infra,
Visualizations, Dashboards

COMPOSER



Programming Language InterOp
Integrate & Manage different code
(C++, Python, Scala, TF, PyT, Spark, ...)
all into the same ML app

FUGUE

```
doubler.py
1 import fugue.processor
2
3 @fugue.processor.register(declaration="doubler.yml")
4 def double(data, models, conf, column_selectors):
5     input_gen = data['Input Dataset']['Data']
6     selected_column = column_selectors['Input Dataset']['to_double'][0]
7
8     def gen():
9         while True:
10            dataset, iter_info = next(input_gen)
11            output_column_name = f'{selected_column}_doubled'
12            dataset[output_column_name] = [float(x * 2) for x in dataset[selected_column]]
13            yield dataset
14
15     return {'Datasets': {'Output Dataset': gen()}}
```

Infra Management
Authentication/Security
Containers
High-Availability
Distributed & Cloud Storage
Cloud & On-Premise Deployments

CHIRON



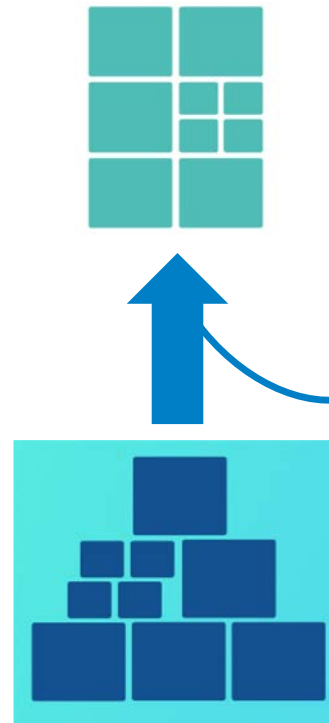
AI (Civil) Engineering with Petuum

Industry Agnostic



 Robo-radiologist	 Insurance Auto-Report	 Virtual EA
 Smart Expense Reports	 Smart Catalog	 Robot Store Staff

Building AI Like Lego

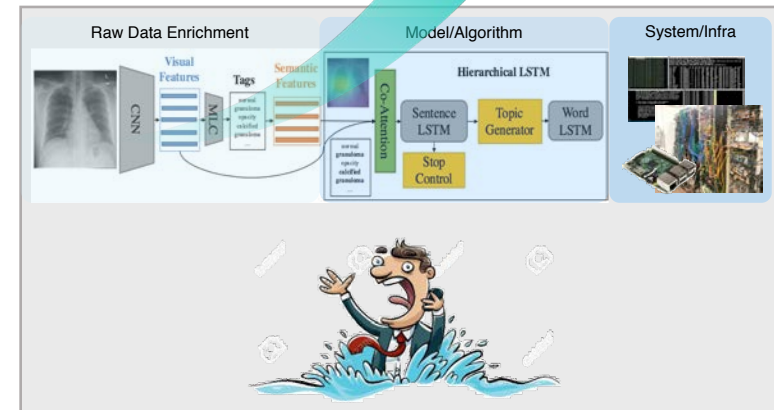
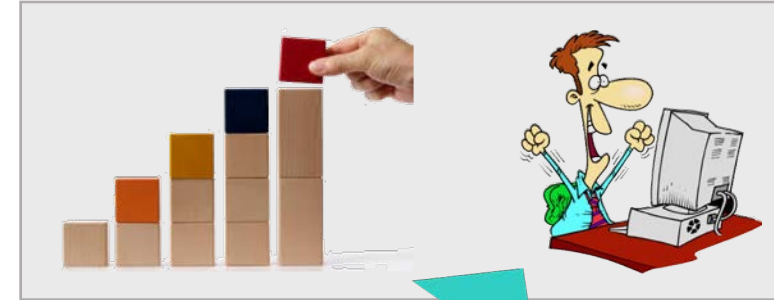


Completed software

10% White-Glove Assembly

90% Completed building blocks

AI With No Tears



Petuum OS

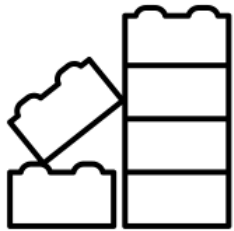
Petuum's Mission

Industrialize AI technology

– turning it from black-box artisanship into standardized engineering process

Transform enterprises across industries

– turning them into owners, builders, and informed users of AI



Sustainable &
Standardized Building
Blocks



One foundation for your
current and future AI-
building needs



A 2018 WEF
Technology Pioneer
Winner



WE ARE HIRING:

ML engineer/manager

Software engineer/manager

System engineer

UI engineer

...

Thank You