# Exemplar Auditing Lifecycle



Static Model

Train

Introspect & update database

Deploy

Build exemplar database

Exemplar Database

Eval on unseen, held-out test

Update database

Introspect dev set (model & data)
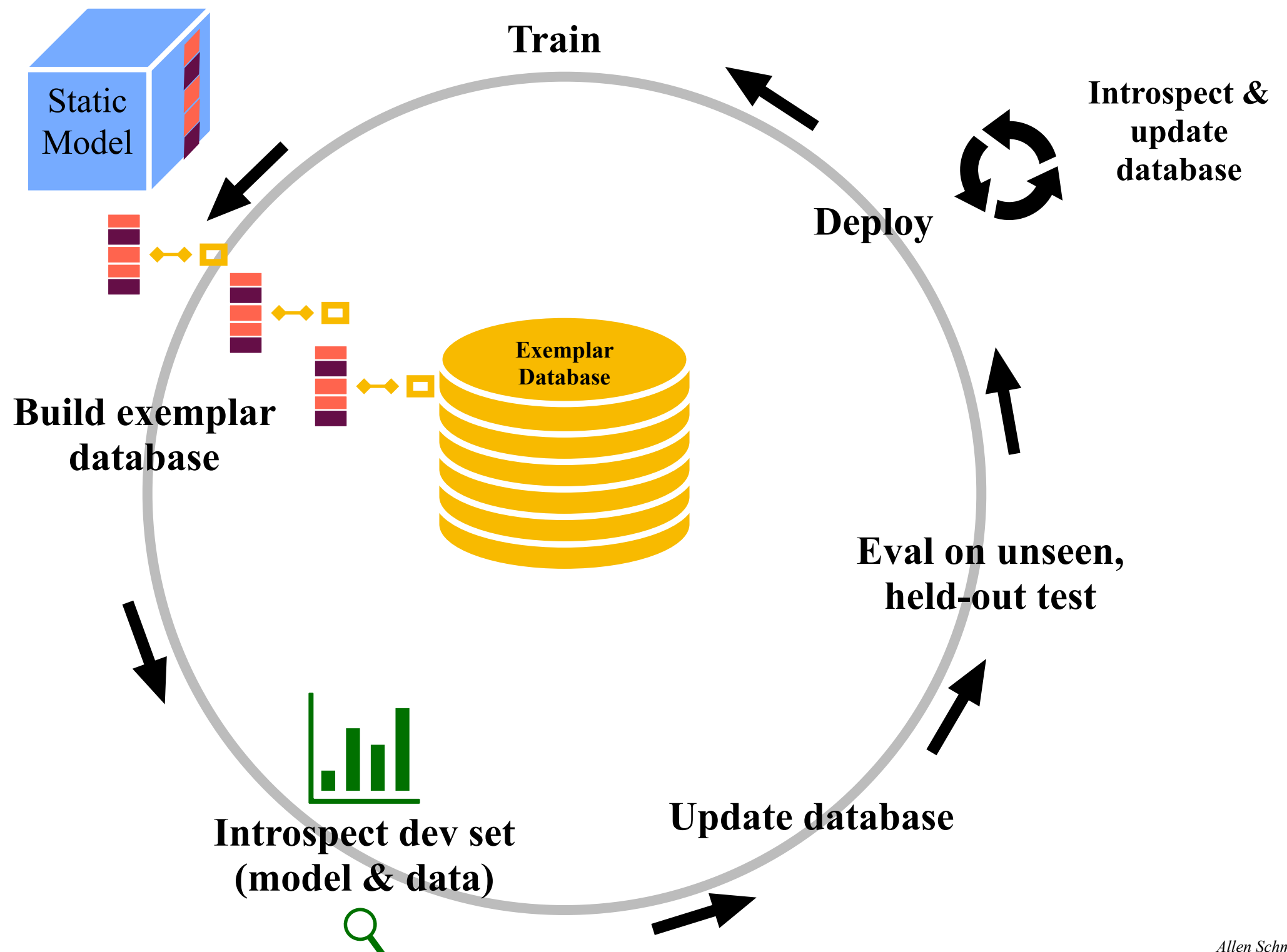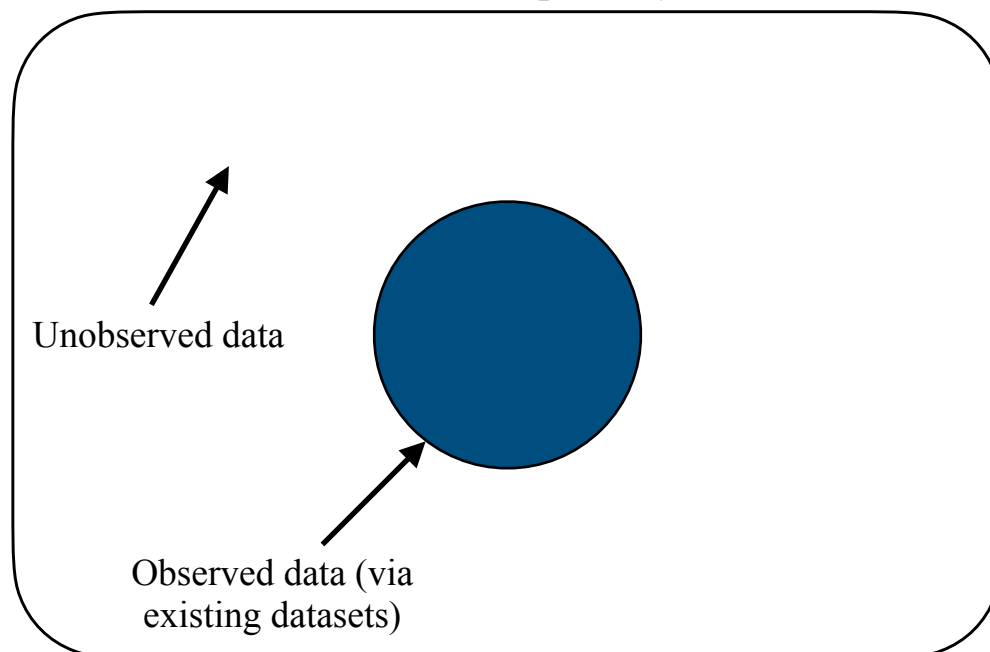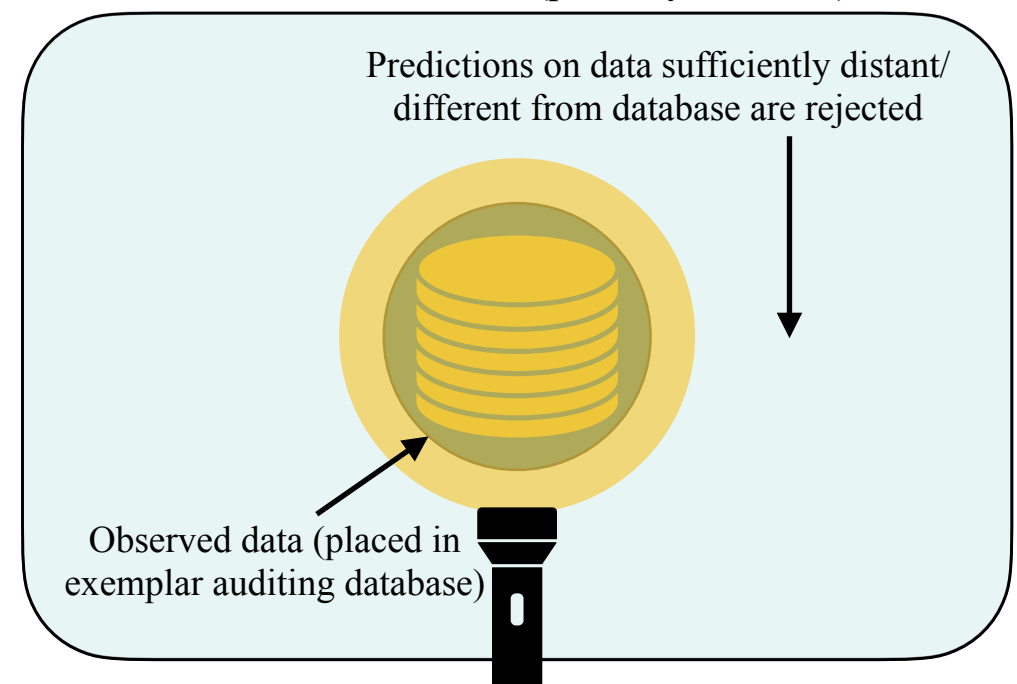
# Out-of-Domain Settings

- Pre-train with as much data as possible

- Add as much data as possible to the database

  - Corral the in-domain space, around the ball of the observed data

  - Never predict over out-of-domain data in high-risk settings. Instead: Rearrange the deployment to handle non-admitted predictions.

**Data distribution for task (partially observed)**
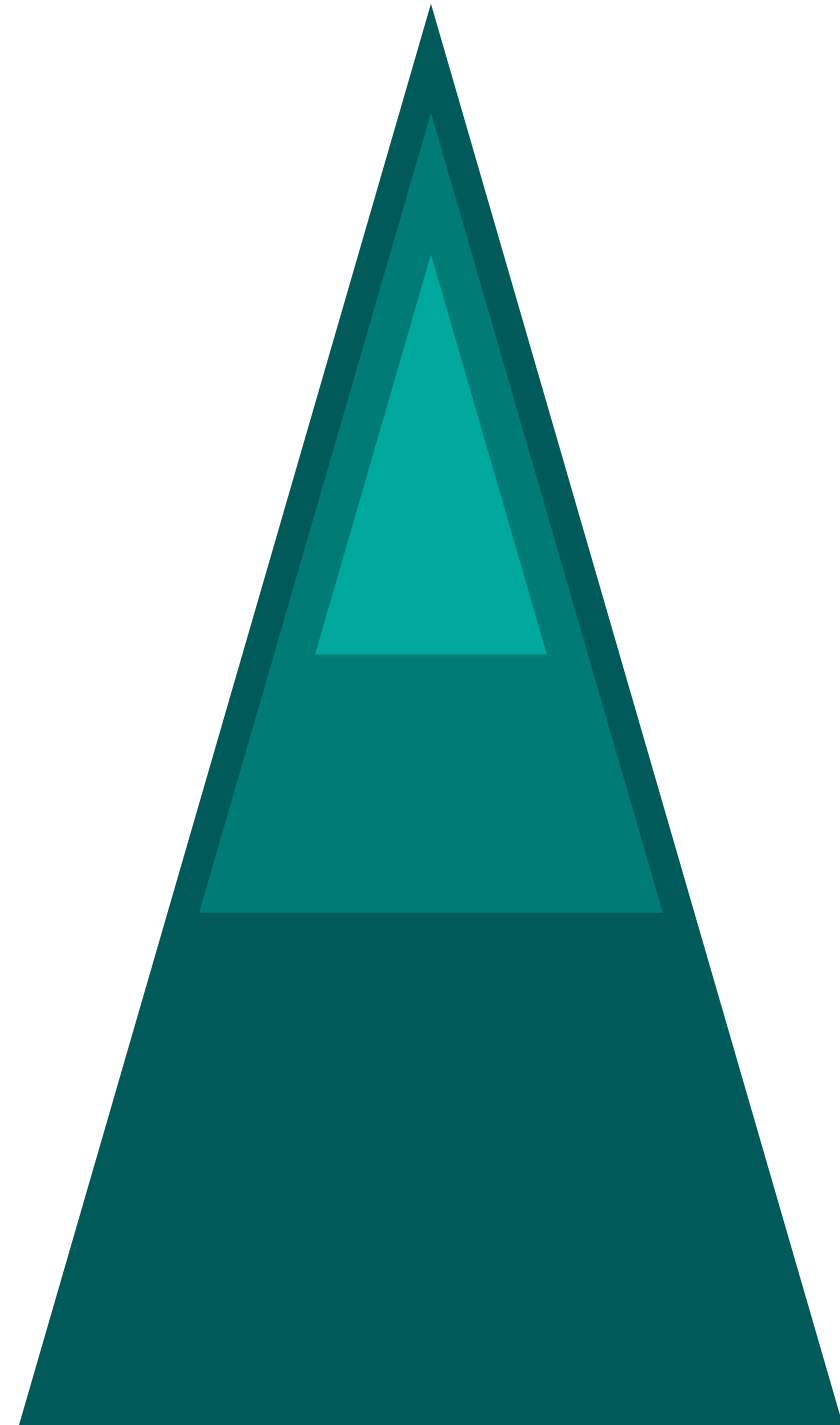
Unobserved data

Observed data (via existing datasets)

**Data distribution for task (partially observed)**

Predictions on data sufficiently distant/ different from database are rejected

Observed data (placed in exemplar auditing database)

# Implementations

- Binary classification: $f : X \rightarrow \{0,1\}$

    - "Detecting Local Insights from Global Labels: Supervised & Zero-Shot Sequence Labeling via a Convolutional Decomposition"

- Multi-label classification: $f : X \rightarrow 2^{|Y|}$

    - "Exemplar Auditing for Multi-Label Biomedical Text Classification"

- Retrieval-classification: $f : X \times \mathcal{D} \rightarrow \left\langle \{0,1,2\}, 2^{|D|} \right\rangle$

    - "Coarse-to-Fine Memory Matching for Joint Retrieval and Classification"

# An End-to-End Retrieval-Classification Model via a Coarse-to-Fine Search over Dense Representations

$min\ L^2\ distance$

Initial, coarse bi-encoder search

CNN

CNN

Shared Transformer LM

Shared Transformer LM

**Seq Q**

$\forall\ S :$   **Seq S**

$min\ L^2\ distance$

CNN

CNN

Subsequent search levels employ cross-encoder search, over the winnowed sequences from earlier levels

Shared Transformer LM

Shared Transformer LM

**Seq Q**   $\forall\ S$ **from previous level** $\Big\}$ **concat(Seq Q, Seq S)**

# Joint Retrieval and Classification Training



Minimize/maximize difference to correct/incorrect matches

$$\boldsymbol{\delta}_L = \left| \boldsymbol{g}^q - \boldsymbol{g}^s \right| \in \mathbb{R}^M$$

Iterative freezing

CNN

CNN

Shared Transformer LM

Shared Transformer LM

Backprop through all search levels

Seq Q

Seq S

The training set is dynamically created via coarse-to-fine search to find hard negatives, as well as prediction sequences that emulate inference

Yields a single model for both retrieval and classification

# Multi-Sequence Representation Composition for Exemplar Auditing

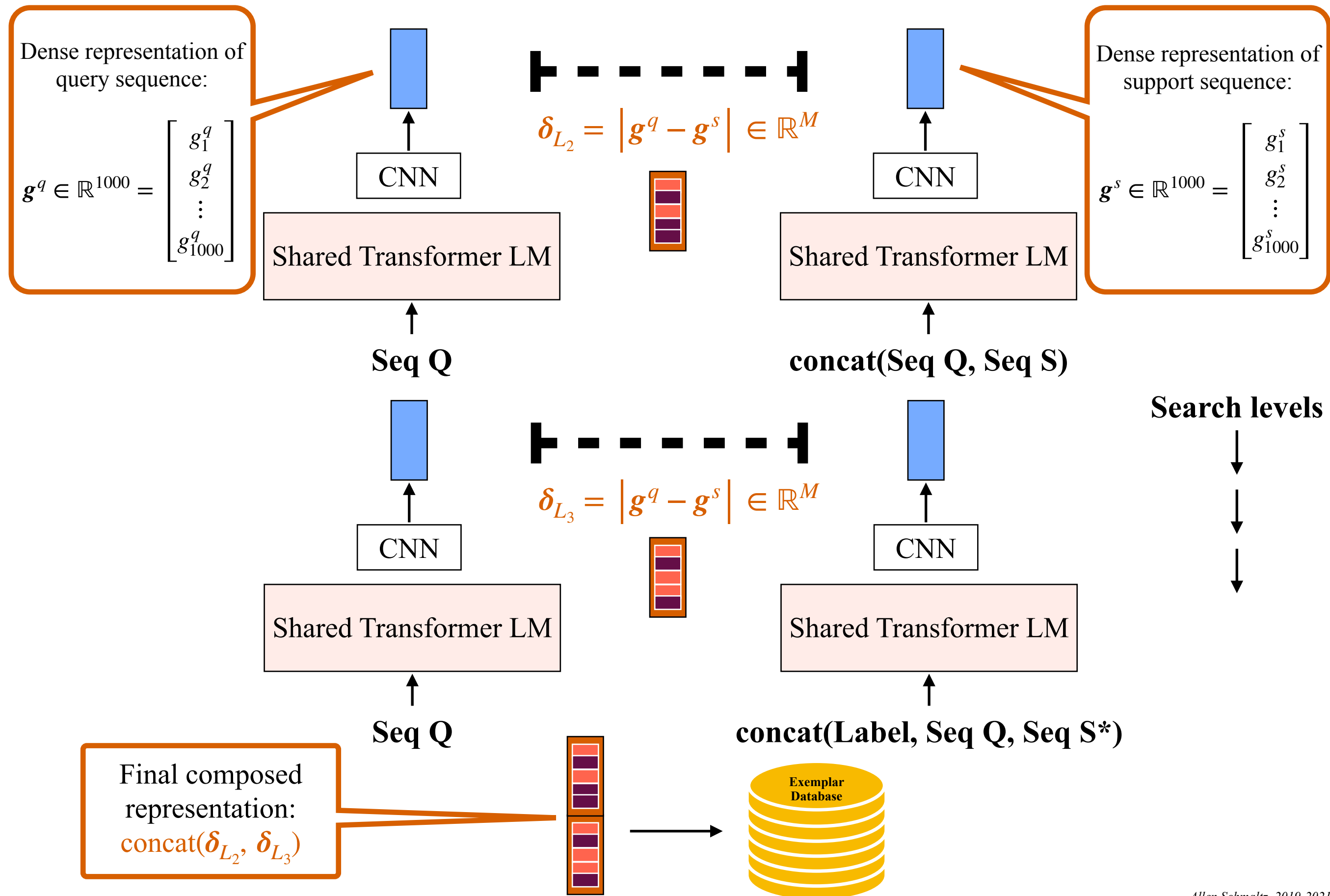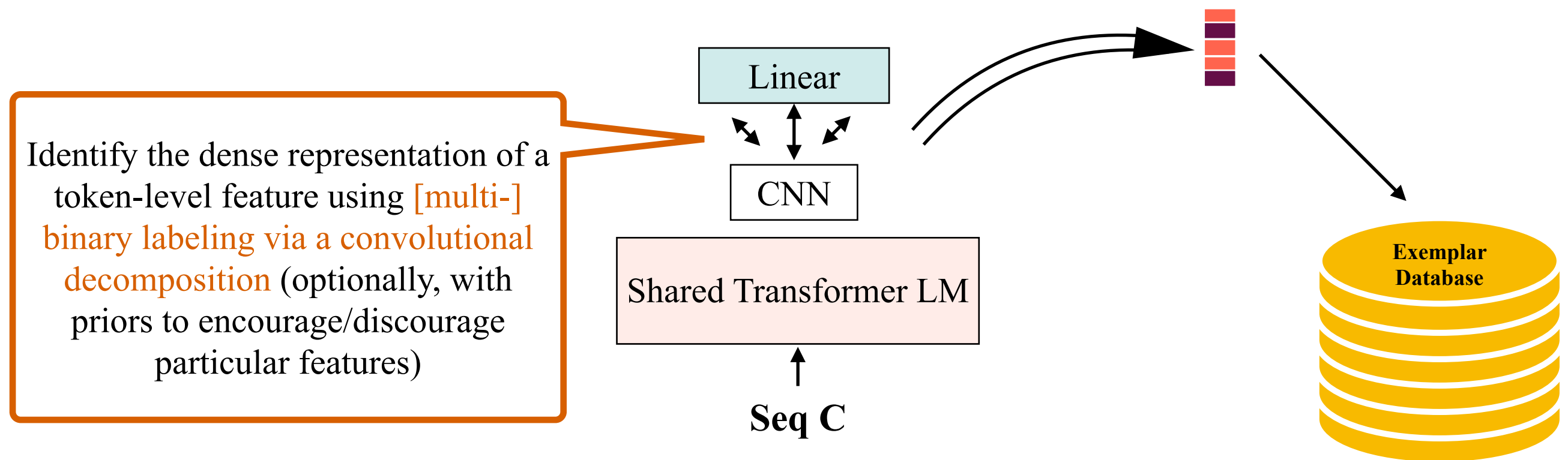**Dense representation of query sequence:**

$$\boldsymbol{g}^q \in \mathbb{R}^{1000} = \begin{bmatrix} g_1^q \\ g_2^q \\ \vdots \\ g_{1000}^q \end{bmatrix}$$

$$\boldsymbol{\delta}_{L_2} = \left| \boldsymbol{g}^q - \boldsymbol{g}^s \right| \in \mathbb{R}^M$$

**Dense representation of support sequence:**

$$\boldsymbol{g}^s \in \mathbb{R}^{1000} = \begin{bmatrix} g_1^s \\ g_2^s \\ \vdots \\ g_{1000}^s \end{bmatrix}$$

CNN

Shared Transformer LM

**Seq Q**

CNN

Shared Transformer LM

**concat(Seq Q, Seq S)**

**Search levels**

$$\boldsymbol{\delta}_{L_3} = \left| \boldsymbol{g}^q - \boldsymbol{g}^s \right| \in \mathbb{R}^M$$

CNN

Shared Transformer LM

**Seq Q**

CNN

Shared Transformer LM

**concat(Label, Seq Q, Seq S*)**

**Final composed representation:**

$$\mathrm{concat}(\boldsymbol{\delta}_{L_2}, \boldsymbol{\delta}_{L_3})$$

**Exemplar Database**

*Allen Schmaltz, 2019-2021*

# Token-Level Representations for Exemplar Auditing

Identify the dense representation of a token-level feature using [multi-] binary labeling via a convolutional decomposition (optionally, with priors to encourage/discourage particular features)

Linear

CNN

Shared Transformer LM

**Seq C**

**Exemplar Database**

*Prospective Outlook*: Interlocking distance constraints across input modalities and tasks via a single, shared model and a dense database…

Productive multi-task outlook, since we get practical models & data analyses along the way

Myoglobin (image from Wikipedia)

*Allen Schmaltz, 2019-2021*