# Matrix-Level Documentation of `gline-rs` Processing Steps

Frédérik Bilhaut

This documents aims at providing a matrix-level description of the pipeline needed for GLiNER inferences, as implemented by `gline-rs`.

Concrete examples are provided for each step, all of which build on the input given in the first one.

## Contents

# 1. Pre-Processing (Common)

## 1.1. Text Input

This is the "user" input for the whole processing.

### 1.1.1. Source Code

- Related struct: `gliner::model::input::text::TextInput`

### 1.1.2. Format

- $n$: number of input texts
- $k$: number of entity class labels
- $I$: sequence of input texts matrix of type `string` and size $n$
- $E$: entity class labels matrix, of type `string` and size $k$

$$I = \begin{bmatrix} \text{text}_1 \\ \text{text}_2 \\ \vdots \\ \text{text}_n \end{bmatrix}$$

$$E = \begin{bmatrix} \text{label}_1 \\ \text{label}_2 \\ \vdots \\ \text{label}_k \end{bmatrix}$$

### 1.1.3. Example

$$I = \begin{bmatrix} \text{"My name is James Bond"} \\ \text{"I like to drive my Aston Martin"} \end{bmatrix}$$

$$E = \begin{bmatrix} \text{"movie character"} \\ \text{"vehicle"} \end{bmatrix}$$

## 1.2. Word-Level Tokenization

### 1.2.1. Transformation

$(I, E) \rightarrow (T, E)$

### 1.2.2. Source Code

- Struct: `gliner::model::input::tokenized::TokenizedInput`
- Transformation: `gliner::model::input::prompt::RawToTokenized`

### 1.2.3. Format

- $n, k$: same as before
- $T$: sequence of sequence of tokenized input texts, of type `string` and size $n$
- $E$: same as before

$$T = \begin{bmatrix} \begin{bmatrix} \text{token}_{1,1} & \text{token}_{1,2} & ... \end{bmatrix} \\ \begin{bmatrix} \text{token}_{2,1} & \text{token}_{2,2} & ... \end{bmatrix} \\ \vdots \\ \begin{bmatrix} \text{token}_{n,1} & \text{token}_{n,2} & ... \end{bmatrix} \end{bmatrix}$$

### 1.2.4. Example

$$T = \begin{bmatrix} \begin{bmatrix} \text{"My"} & \text{"name"} & \text{"is"} & \text{"James"} & \text{"Bond"} \end{bmatrix} \\ \begin{bmatrix} \text{"I"} & \text{"like"} & \text{"to"} & \text{"drive"} & \text{"my"} & \text{"Aston"} & \text{"Martin"} \end{bmatrix} \end{bmatrix}$$

### 1.3. Prompt Preparation

Prepared prompts, appending entity and text tokens.

### 1.3.1. Transformation

$(T, E) \rightarrow P$

### 1.3.2. Source Code

- Struct: `gliner::model::input::prompt::PromptInput`
- Transformation from `TokenizedInput`: `gliner::model::input::prompt::TokenizedToPrompt`

### 1.3.3. Format

$$P = \begin{bmatrix} \left[ \text{<<ENT>>} \quad \text{label}_{1,1} \quad \text{<<ENT>>} \quad \text{label}_{1,2} \quad \dots \quad \text{<<SEP>>} \quad \text{token}_{1,1} \quad \text{token}_{1,2} \quad \dots \right] \\ \left[ \text{<<ENT>>} \quad \text{label}_{2,1} \quad \text{<<ENT>>} \quad \text{label}_{2,2} \quad \dots \quad \text{<<SEP>>} \quad \text{token}_{2,1} \quad \text{token}_{2,2} \quad \dots \right] \\ \vdots \\ \left[ \text{<<ENT>>} \quad \text{label}_{n,1} \quad \text{<<ENT>>} \quad \text{label}_{n,2} \quad \dots \quad \text{<<SEP>>} \quad \text{token}_{n,1} \quad \text{token}_{n,2} \quad \dots \right] \end{bmatrix}$$

### 1.3.4. Example

$$P = \begin{bmatrix} \left[ \text{<<ENT>>} \quad \text{"movie character"} \quad \text{<<ENT>>} \quad \text{"vehicle"} \quad \dots \quad \text{<<SEP>>} \quad \text{"My"} \quad \text{"name"} \quad \text{"is"} \quad \text{"James"} \quad \text{"Bond"} \right] \\ \left[ \text{<<ENT>>} \quad \text{"movie character"} \quad \text{<<ENT>>} \quad \text{"vehicle"} \quad \dots \quad \text{<<SEP>>} \quad \text{"I"} \quad \text{"like"} \quad \text{"to"} \quad \text{"drive"} \quad \text{"my"} \quad \text{"Austin"} \quad \text{"Martin"} \right] \end{bmatrix}$$

### 1.4. Prompt Encoding (Sub-Word Tokenization)

#### 1.4.1. Transformation
$P \rightarrow (I, A, W, L)$

#### 1.4.2. Source Code
- Struct: `gliner::model::input::encoded::EncodedPrompt`
- Transformation: `gliner::model::input::encoded::PromptsToEncoded`

#### 1.4.3. Format
- k: maximum number of sub-word tokens within a sequence, adding start (1) and end (2) tokens
- I: encoded prompts of type `i64` and shape $(n * k)$
- A: attention masks of type `i64` and shape $(n * k)$
- W: word masks of type `i64` and shape $(n * k)$
- L: text lengths of type `i64` and shape $(n * 1)$

$$I = \begin{pmatrix} \text{token\_id}_{1,1} & \text{token\_id}_{1,2} & \dots & \text{token\_id}_{1,k} \\ \text{token\_id}_{2,1} & \text{token\_id}_{2,2} & \dots & \text{token\_id}_{2,k} \\ \vdots & \vdots & \ddots & \vdots \\ \text{token\_id}_{n,1} & \text{token\_id}_{n,2} & \dots & \text{token\_id}_{n,k} \end{pmatrix}$$

$$A = \begin{pmatrix} \text{attn\_mask}_{1,1} & \text{attn\_mask}_{1,2} & \dots & \text{attn\_mask}_{1,k} \\ \text{attn\_mask}_{2,1} & \text{attn\_mask}_{2,2} & \dots & \text{attn\_mask}_{2,k} \\ \vdots & \vdots & \ddots & \vdots \\ \text{attn\_mask}_{n,1} & \text{attn\_mask}_{n,2} & \dots & \text{attn\_mask}_{n,k} \end{pmatrix}$$

$$W = \begin{pmatrix} \text{word\_mask}_{1,1} & \text{word\_mask}_{1,2} & \dots & \text{word\_mask}_{1,k} \\ \text{word\_mask}_{2,1} & \text{word\_mask}_{2,2} & \dots & \text{word\_mask}_{2,k} \\ \vdots & \vdots & \ddots & \vdots \\ \text{word\_mask}_{n,1} & \text{word\_mask}_{n,2} & \dots & \text{word\_mask}_{n,k} \end{pmatrix}$$

$$L = \begin{pmatrix} l_1 \\ \vdots \\ l_n \end{pmatrix}$$

#### 1.4.4. Example

$$I = \begin{pmatrix} 1 & 128002 & 1421 & 1470 & 128002 & 1508 & 128003 & 573 & 601 & 269 & 1749 & 8728 & 2 & 0 & 0 \\ 1 & 128002 & 1421 & 1470 & 128002 & 1508 & 128003 & 273 & 334 & 264 & 1168 & 312 & 20844 & 2963 & 2 \end{pmatrix}$$

$$A = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix}$$

$$W = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 2 & 3 & 4 & 5 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 0 \end{pmatrix}$$

$$L = \begin{pmatrix} 5 \\ 7 \end{pmatrix}$$

# 2. Pre-Processing (Span Mode)
Downstream of the aforementioned steps.

## 2.1. Span Preparation

### 2.1.1. Transformation
$$(I, A, W, L) \rightarrow (I, A, W, L, S_I, S_M)$$

### 2.1.2. Format
- $n, k, I, A, W, L$: same as before.
- $s$: maximum possible number of spans for one sequence
- $S_I$: span offsets, of type `i64` and shape $(n * s * 2)$
- $S_M$: span masks, of type `bool` and shape $(n * s)$

$$S_I = \begin{pmatrix} (\text{start}_{1,1} \; \text{end}_{1,1}) & (\text{start}_{1,2} \; \text{end}_{1,2}) & \dots & (\text{start}_{1,s} \; \text{end}_{1,s}) \\ (\text{start}_{2,1} \; \text{end}_{2,1}) & (\text{start}_{2,2} \; \text{end}_{2,2}) & \dots & (\text{start}_{2,s} \; \text{end}_{2,s}) \\ \vdots & \vdots & \ddots & \vdots \\ (\text{start}_{n,1} \; \text{end}_{n,1}) & (\text{start}_{n,2} \; \text{end}_{n,2}) & \dots & (\text{start}_{n,s} \; \text{end}_{n,s}) \end{pmatrix}$$

$$S_M = \begin{pmatrix} \text{span\_mask}_{1,1} & \text{span\_mask}_{1,2} & \dots & \text{span\_mask}_{1,s} \\ \text{span\_mask}_{2,1} & \text{span\_mask}_{2,2} & \dots & \text{span\_mask}_{2,s} \\ \vdots & \vdots & \ddots & \vdots \\ \text{span\_mask}_{n,1} & \text{span\_mask}_{n,2} & \dots & \text{span\_mask}_{n,s} \end{pmatrix}$$

### 2.1.3. Example
Note: for readability purposes, inside matrices are split into rows (one per token) but they are actually in one dimension $s$ (see format above).

$$S_I = \begin{pmatrix} \begin{pmatrix} (0\ 0) & (0\ 1) & (0\ 2) & (0\ 3) & (0\ 4) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & \hookleftarrow \\ (1\ 1) & (1\ 2) & (1\ 3) & (1\ 4) & (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \hookleftarrow \\ (2\ 2) & (2\ 3) & (2\ 4) & (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \hookleftarrow \\ (3\ 3) & (3\ 4) & (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \hookleftarrow \\ (4\ 4) & (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \hookleftarrow \\ (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \hookleftarrow \\ (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \end{pmatrix} \\ \begin{pmatrix} (0\ 0) & (0\ 1) & (0\ 2) & (0\ 3) & (0\ 4) & (0\ 5) & (0\ 6) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & \hookleftarrow \\ (1\ 1) & (1\ 2) & (1\ 3) & (1\ 4) & (1\ 5) & (1\ 6) & (0\ 0) & \dots & \dots & \dots & \dots & \dots & \hookleftarrow \\ (2\ 2) & (2\ 3) & (2\ 4) & (2\ 5) & (2\ 6) & (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \hookleftarrow \\ (3\ 3) & (3\ 4) & (3\ 5) & (3\ 6) & (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \hookleftarrow \\ (4\ 4) & (4\ 5) & (4\ 6) & (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \hookleftarrow \\ (5\ 5) & (5\ 6) & (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \hookleftarrow \\ (6\ 6) & (0\ 0) & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \end{pmatrix} \end{pmatrix}$$

$$S_M = \left( \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \atop \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \hookleftarrow \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \right)$$

## 3. Pre-Processing (Token Mode)

Nothing to be done beside the common steps.

# 4. Post-Processing (Span Mode)

## 4.1. Logits Output

### 4.1.1. Source Code

- Struct: `gliner::model::output::TensorOutput`

### 4.1.2. Format

- $n$: number of text sequences
- $w$: maximum number of tokens in one sequence
- $s$: maximum number of possible spans for one token (seee above)
- $k$: number of entity labels
- $O$: logits output, of type `f32` and shape $(n * w * s * k)$
- $v_{n,w,s,k}$: raw model output for sequence $n$, token $w$, span $s$ and label $k$.

$$O = \begin{pmatrix} \left( \begin{pmatrix} v_{1,1,1,1} & \cdots & v_{1,1,1,k} \\ & \vdots & \\ v_{1,1,s,1} & \cdots & v_{1,1,s,k} \end{pmatrix} \cdots \begin{pmatrix} v_{1,w,1,1} & \cdots & v_{1,w,1,k} \\ & \vdots & \\ v_{1,w,s,1} & \cdots & v_{1,w,s,k} \end{pmatrix} \right) \\ \vdots \\ \left( \begin{pmatrix} v_{n,1,1,1} & \cdots & v_{n,1,1,k} \\ & \vdots & \\ v_{n,1,s,1} & \cdots & v_{n,1,s,k} \end{pmatrix} \cdots \begin{pmatrix} v_{n,w,1,1} & \cdots & v_{n,w,1,k} \\ & \vdots & \\ v_{n,w,s,1} & \cdots & v_{n,w,s,k} \end{pmatrix} \right) \end{pmatrix}$$

### 4.1.3. Example

In this case $s = 12$. For readability purposes, the raw values are "sigmoided" ($S(x) = \frac{1}{1+e^{-x}}$) and then "ReLUed" with a threshold $t = 0.5$.

$$O_{S,t} = \begin{pmatrix} \begin{pmatrix} (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0.89\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \end{pmatrix} \\ \begin{pmatrix} (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0.96) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \\ (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) & (0\ 0) \end{pmatrix} \end{pmatrix}$$

Which means:

- In the 1st sequence, the span starting with the 4th token and ending with the 5th one has a probability of 0.89 to match the 1st entity class.
- In the 2nd sequence, the span starting with the 6th token and ending with the 7th one has a probability of 0.96 to match the second 2nd class.

## 4.2. Span Decoding

### 4.2.1. Transformation

$(O, L) \rightarrow S$

### 4.2.2. Source Code

- Struct: `gliner::model::output::decoded::SpanOutput`
- Transformation: `gliner::model::output::decoded::span::TensorsToDecoded`

### 4.2.3. Format

- $t$: threshold
- $n$: number of input sequences
- $L$: text lengths as defined before
- $S$: sequence of spans $(i, j, k, p)$ where:
  - $i$ is the index of the first token of sequence $m$ with $i < j$ and $i < L(m)$
  - $j$ is the index of the last token with the same constraints as $i$
  - $k$ is the entity class,
  - $p$ is the probability for class $k$ with $p \geq t$

$$S = \begin{bmatrix} \left[ (i_{1,1}, j_{1,1}, k_{1,1}, p_{1,1}) \quad (i_{1,2}, j_{1,2}, k_{1,2}, p_{1,2}) \quad \ldots \right] \\ \vdots \\ \left[ (i_{n,1}, j_{n,1}, k_{n,1}, p_{n,1}) \quad (i_{n,2}, j_{n,2}, k_{n,2}, p_{n,2}) \quad \ldots \right] \end{bmatrix}$$

### 4.2.4. Example

$$S = \begin{bmatrix} \left[ (4,5,1,0.89) \right] \\ \left[ (6,7,2,0.96) \right] \end{bmatrix}$$

# 5. Post-Processing (Token Mode)

## 5.1. Logits Output

### 5.1.1. Source Code
- Struct: `gliner::model::output::TensorOutput`

### 5.1.2. Format
- $n$: number of text sequences
- $w$: maximum number of tokens in one sequence
- $k$: number of entity labels
- $O$: logits output, of type `f32` and shape $(3 * n * w * k)$ with:
  - $s_{n,w,k}$: raw model output for a start token $w$ in sequence $n$ and label $k$.
  - $e_{n,w,k}$: raw model output for an end token $w$ in sequence $n$ and label $k$.
  - $i_{n,w,k}$: raw model output for an inside token $w$ in sequence $n$ and label $k$.

$$
O = \begin{pmatrix}
\left( \begin{pmatrix} s_{1,1,1} & \cdots & s_{1,1,k} \\ \vdots & \ddots & \vdots \\ s_{1,w,1} & \cdots & s_{1,w,k} \end{pmatrix} \cdots \begin{pmatrix} s_{n,1,1} & \cdots & s_{n,1,k} \\ \vdots & \ddots & \vdots \\ s_{n,w,1} & \cdots & s_{n,w,k} \end{pmatrix} \right) \\
\left( \begin{pmatrix} e_{1,1,1} & \cdots & e_{1,1,k} \\ \vdots & \ddots & \vdots \\ e_{1,w,1} & \cdots & e_{1,w,k} \end{pmatrix} \cdots \begin{pmatrix} e_{n,1,1} & \cdots & e_{n,1,k} \\ \vdots & \ddots & \vdots \\ e_{n,w,1} & \cdots & e_{n,w,k} \end{pmatrix} \right) \\
\left( \begin{pmatrix} i_{1,1,1} & \cdots & i_{1,1,k} \\ \vdots & \ddots & \vdots \\ i_{1,w,1} & \cdots & i_{1,w,k} \end{pmatrix} \cdots \begin{pmatrix} i_{n,1,1} & \cdots & i_{n,1,k} \\ \vdots & \ddots & \vdots \\ i_{n,w,1} & \cdots & i_{n,w,k} \end{pmatrix} \right)
\end{pmatrix}
$$

### 5.1.3. Example
For readability purposes, the raw values are "sigmoided" ($S(x) = \frac{1}{1+e^{-x}}$) and then "ReLUed" with a threshold $t = 0.5$.

$$
O_{S,t} = \begin{pmatrix}
\left( \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0.97 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0.99 \\ 0 & 0 \end{pmatrix} \right) \\
\left( \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0.96 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0.97 \end{pmatrix} \right) \\
\left( \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0.98 & 0 \\ 0.98 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0.99 \\ 0 & 0.99 \end{pmatrix} \right)
\end{pmatrix}
$$

### 5.2. Span Decoding

#### 5.2.1. Transformation
$O \rightarrow S$

#### 5.2.2. Source Code
- Struct: `gliner::model::output::decoded::SpanOutput`
- Transformation: `gliner::model::output::decoded::token::TensorsToDecoded`

#### 5.2.3. Format
Same format as in span-mode.

# 6. Post-Processing (Common)

## 6.1. Span Filtering (Greedy Search)

### 6.1.1. Transformation

$S \rightarrow S'$

### 6.1.2. Source Code

- Struct: `gliner::model::output::decoded::SpanOutput`
- Transformation: `gliner::model::output::decoded::greedy::GreedySearch`

### 6.1.3. Format

Same as span output.