

# メールサービスを面白くするための**AI**と 自然言語処理の基本と応用。

**Futuristic Communication**  
**The Role of AI in Email Services**



---

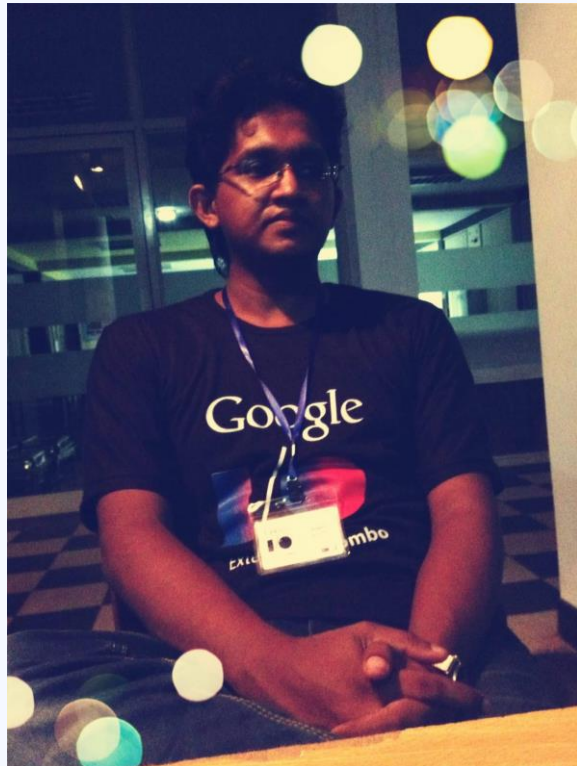
JPAAWG 6TH MEETING - 2023

# Consent.

- ❑ I warmly welcome attendees to share insights and learnings from my presentation at JPAAWG 6th General Meeting on social media. I suggest using the official hashtag when posting on platforms such as Twitter, Facebook, and Instagram to help others easily find and follow the conversation.
- ❑ Photography and video recording are allowed during my presentation, and I appreciate if you could attribute any shared content to me and JPAAWG 6th General Meeting.
- ❑ However, I ask that you respect the privacy and preferences of others in the audience. If someone indicates that they do not wish to be photographed or recorded, please respect their wishes.

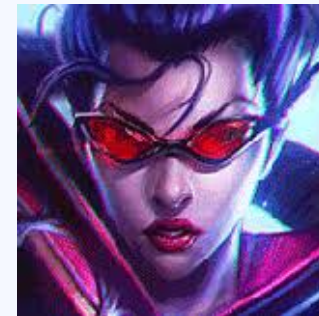
Thank you for your cooperation and understanding.

# About Me



- Name : Nuwan Senevirathne (ヌワン)
- Company : Qualitia
- Position : AI Engineer

- I love gaming
- League of Legend (LOL)
- AD main





# Agenda

**01**

**Neural Networks.**

**02**

**Attention is All You Need.**

**03**

**Subject Line Generation.**

**04**

**Email Mis-sending Prevention.**





# 01

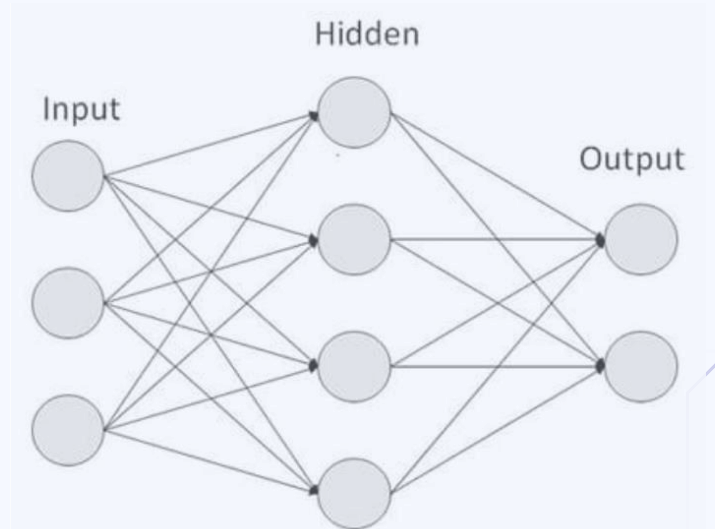
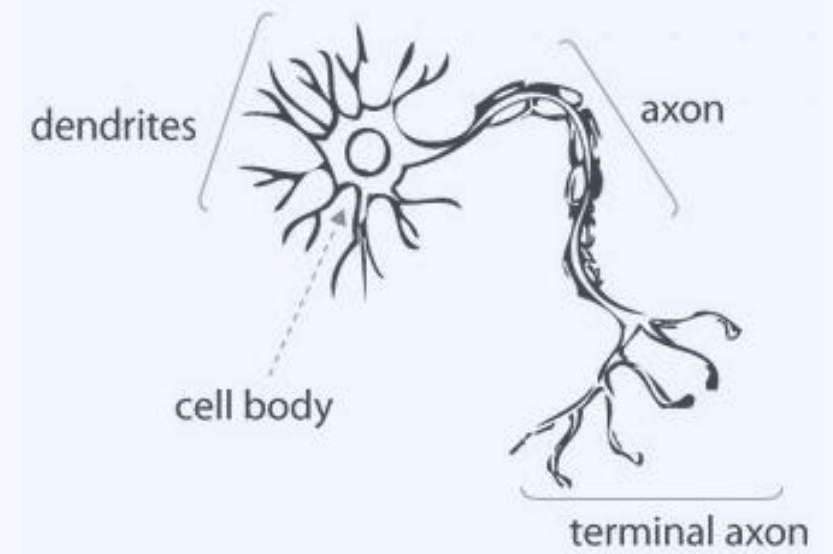
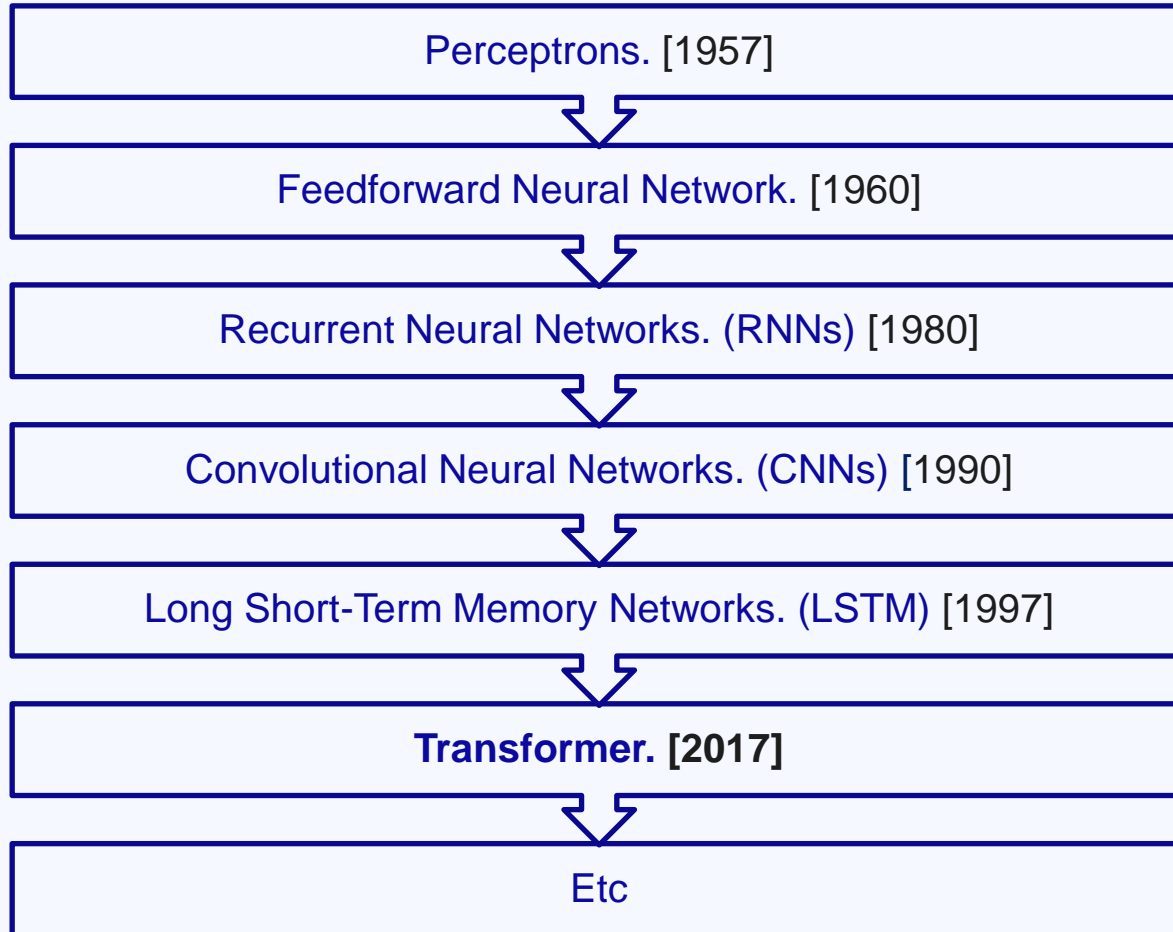
## Neural Networks.

---

**Simple Introduction.**



# Neural Networks.





# 02

## Attention is All You Need.

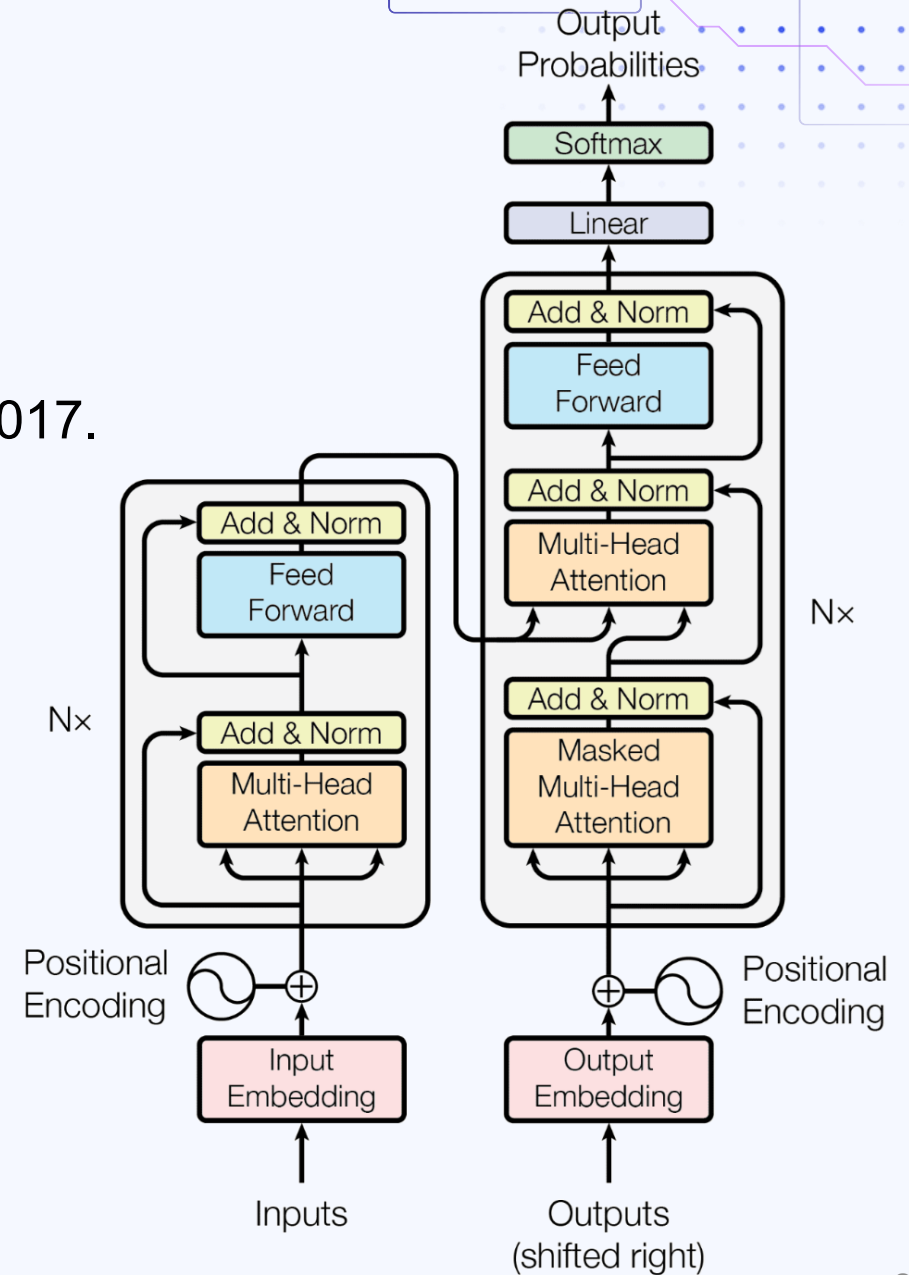
---

### Transformer Model

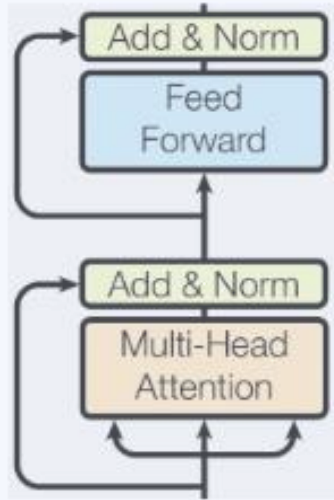


# What is Transformer.

- Neural Network Architecture.
- Vaswani et al. in the paper "Attention is All You Need" 2017.
- Based on attention mechanisms.
- Major component is **multi-head self attention**.
- Useful in;
  - Machine translations.
  - Text Generations.
  - Text Classifications.



# Components of the Transformer Model.



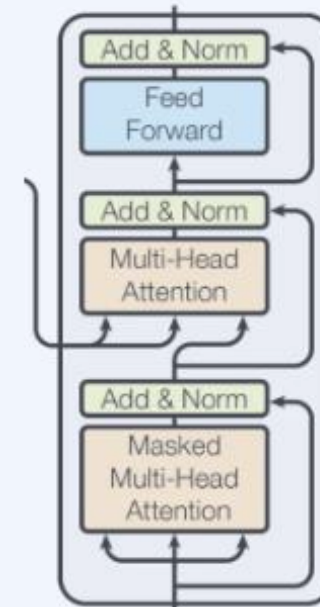
Transformer encoder architecture

## ❑ Encoder

- ❖ Takes Input data.
- ❖ Transform into series of numerical vectors.
- ❖ Each vector capture the meaning of word in the context of the whole sentence.

## ❑ Decoder

- ❖ Takes the vectors.
- ❖ Generate the output.



Transformer decoder architecture

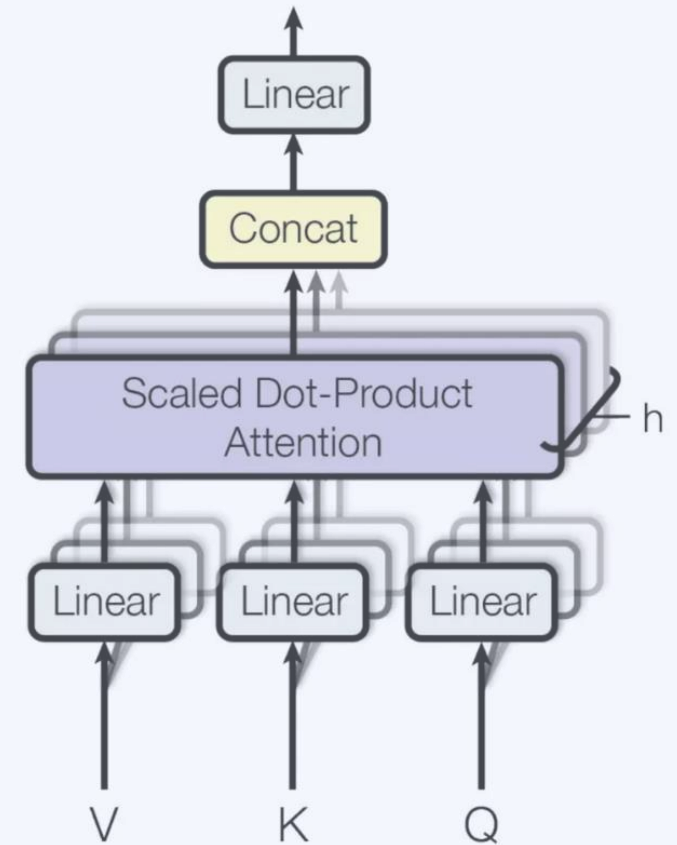
# Attention Mechanism.

## □ General Concept

- Allows a model to focus on specific parts of the input when producing an output.
- Weighs the importance of different parts of the input.
- Decide how much focus to put on different areas.

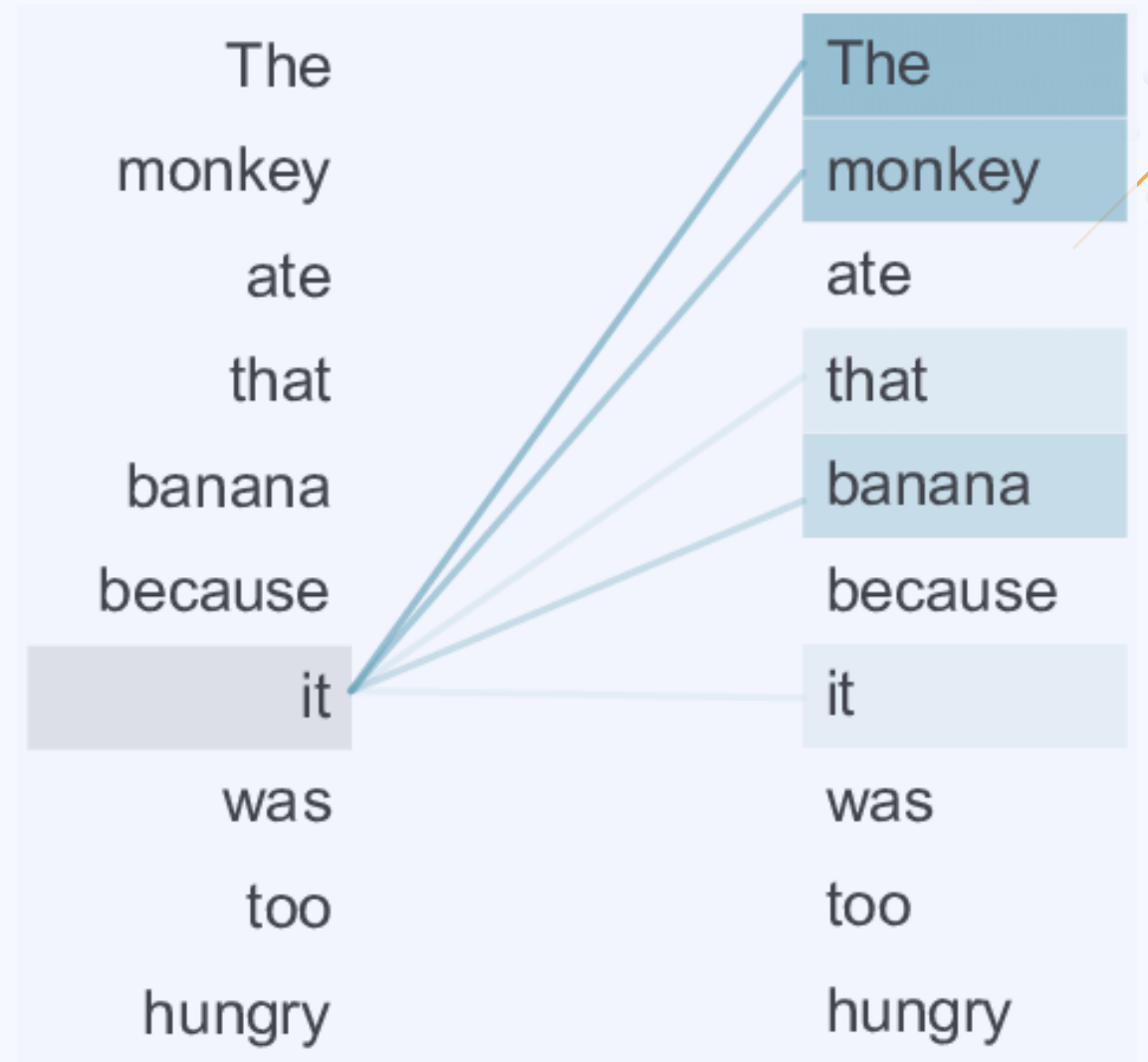
## □ Types

- Scaled dot-product attention,
- Multiplicative attention
- Additive attention
- Etc.



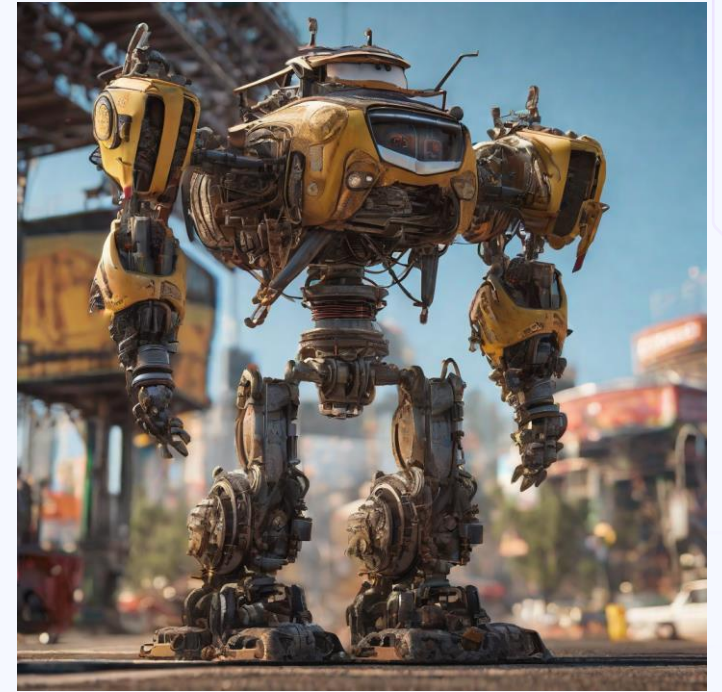
# Self Attention.

- ❑ Scaled dot-product attention.
- ❑ Specific Type of Attention:
  - Model attends to all parts of the input simultaneously
  - Compute the representation of each part in the context of all other parts.
- ❑ Consider the entire context of a sentence.



# Transformer: Conclusion

- ❑ Innovative architecture.
- ❑ Widely used in modern NLP.
- ❑ Understand words
- ❑ Understand context.
- ❑ Understand relationships.
- ❑ Many models followed this model architecture and produce SOTA results on various of NLP tasks.
- ❑ If anyone need more information, visit <https://jalammar.github.io/illustrated-transformer/> . This post explain the things very well.





# AI in Email Services and Security.

- ❑ Spam Detection.
- ❑ Phishing Detection.
- ❑ Malware Detection.
- ❑ Anomaly Detection.
- ❑ Risk Scoring.
- ❑ Data Loss Prevention.
  
- ❑ There are many more things.



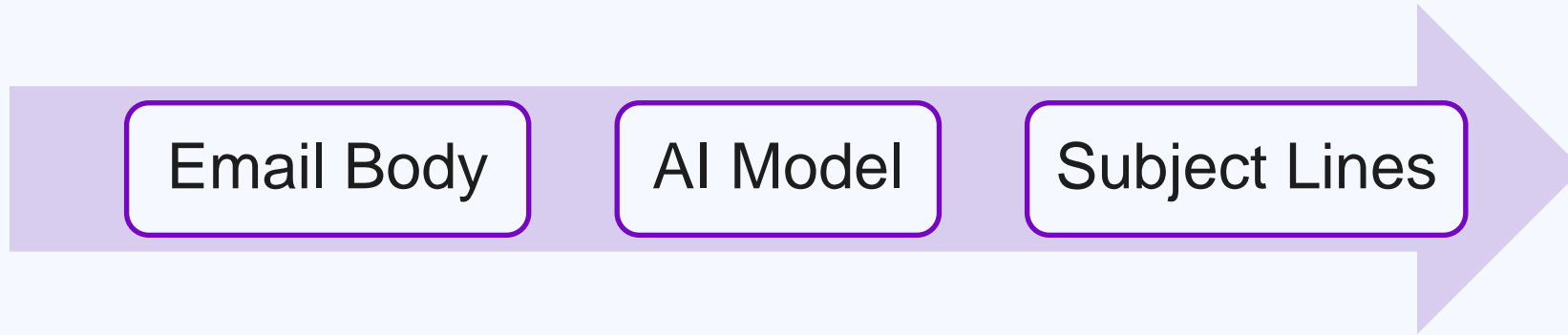
# 03

## Subject Line Generation.

---

**Finetuned custom T5 Model.**

# What is Subject Line Generation?



1. User inputs email content.
2. AI model analyzes the content.
3. AI model generates a list of potential subject lines.
4. More effective approach.
5. Eliminates the need for users to brainstorm a subject line.

# Subject Line Generation Example

各位

お疲れ様です。人事総務部CDです。

年末調整の申告書および添付書類提出期限が  
来週12/3（金）までとなっておりますので  
再度ご案内させていただきます。

既にご提出頂いた方もいらっしゃいますが、  
各自ご対応のほどよろしくお願ひいたします。  
ご不明な点はCDまたはABさんまでお問い合わせください。

※添付書類原本についてはなるべく台紙に貼り付けて頂けると助かります。  
A4サイズの書類はそのまま結構です。

以上、よろしくお願い致します。

Generate

Reset

Output

年末調整の申告書および添付書類提出期限について

年末調整の申告について(再送)

年末調整申告書および添付書類提出期限について

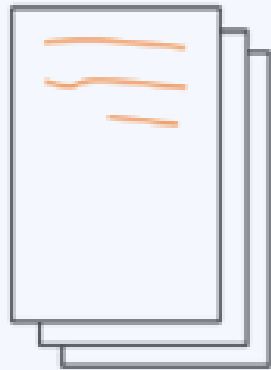
年末調整の申告について(再送) **【重要】**

# Summarization Methods.

Source Document



Extractive Summary



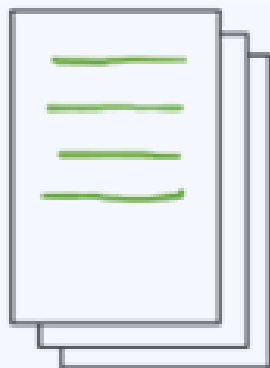
## □ **Extractive summarization.**

- Identify Key phrases or sentences
- Use those extracted phrases as summary.
- Same sentences and structure.

Source Document



Abstractive Summary



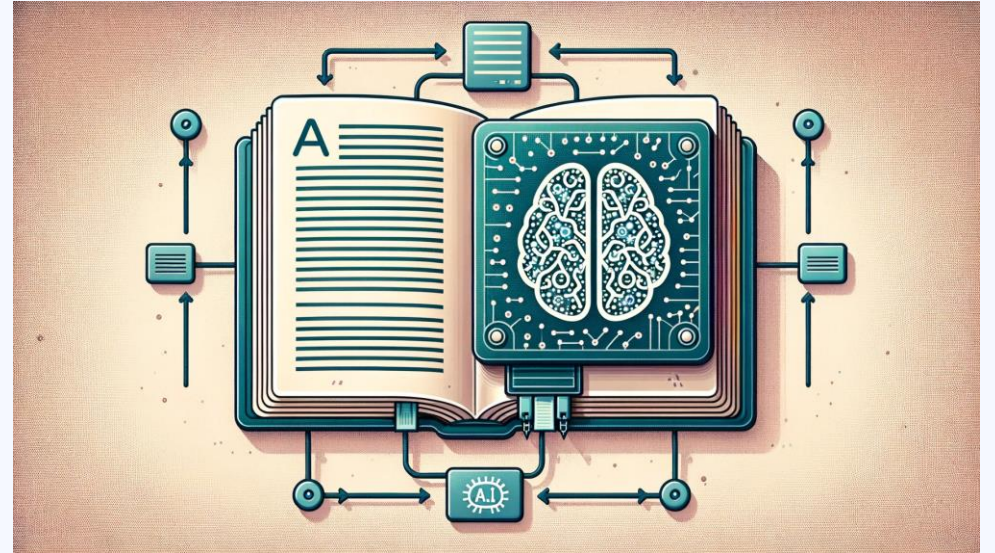
## □ **Abstractive summarization.**

- Paraphrasing.
- Different words, different sentence structure.
- Concise summary.



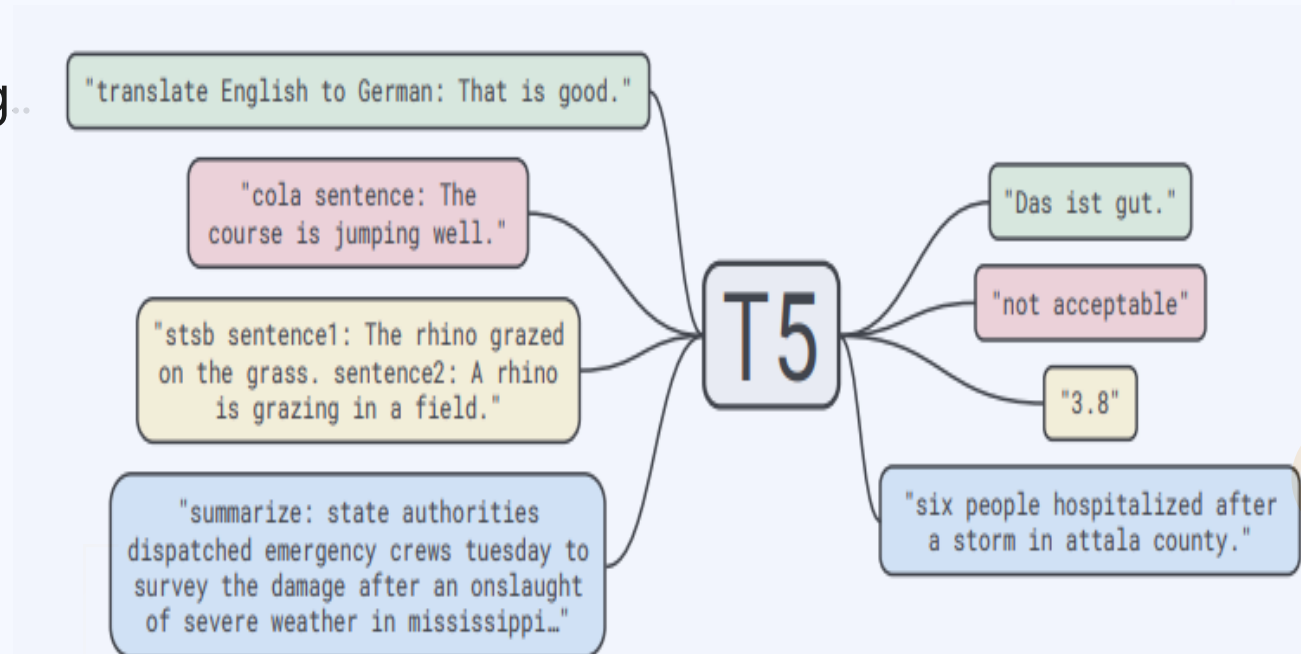
# Why Abstractive Summarization?

- ❖ Short.
  - ❖ Concise.
  - ❖ Directly relevant to the contents.
  - ❖ Capture mainpoint or purpose.
- Example: (Email content)
    - 。 みなさま、お疲れ様です。平野です。
    - 。 **M3AAWG**とは何か、**M3AAWG**でのメールセキュリティや周辺トピックの最新の話、サンフランシスコの最近の雰囲気、などについてお話しします。
  - Subject: "**M3AAWG**概要とサンフランシスコ現地レポートのご案内".

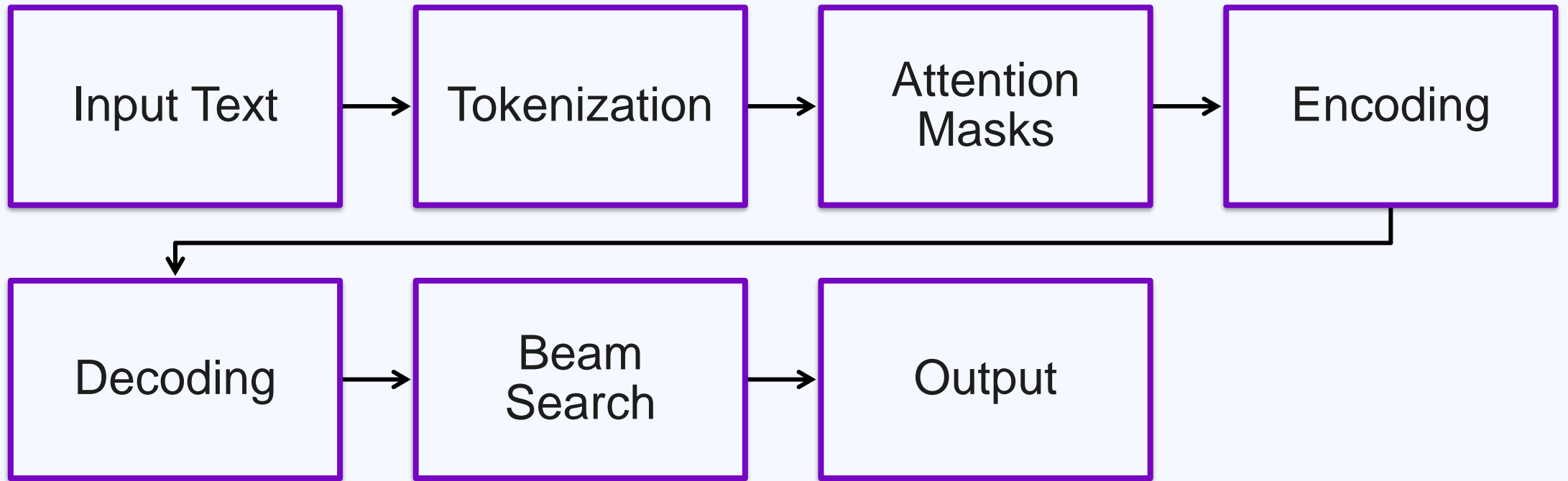


# T5 (Text-To-Text Transfer Transformer) Model.

1. Treats every NLP problem as a text-to-text problem.
2. Causal Language Model.
3. Heart of T5 are self-attention mechanisms.
4. T5 uses Relative Positional Encoding.
5. Causal masking.
6. Encoder-decoder.



# Overview of the T5 process.





## T5: Text embedding flow.

Tokenization

Word Embeddings

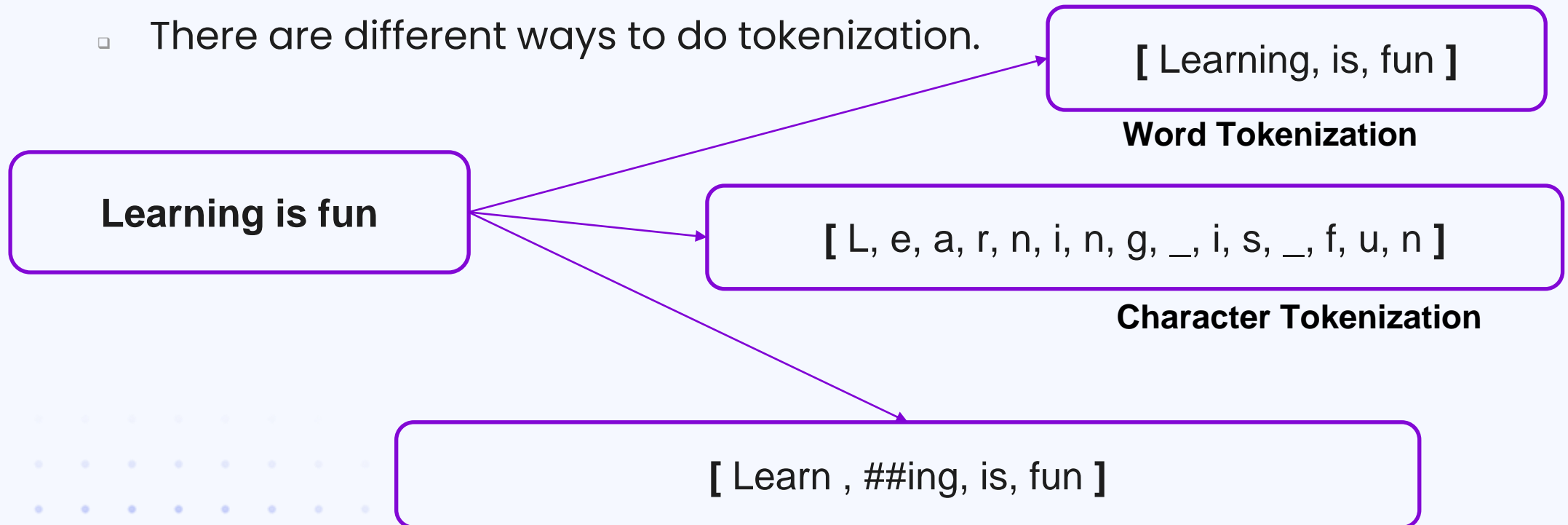
Positional Encodings

Combination of Embeddings

Contextualized Embeddings

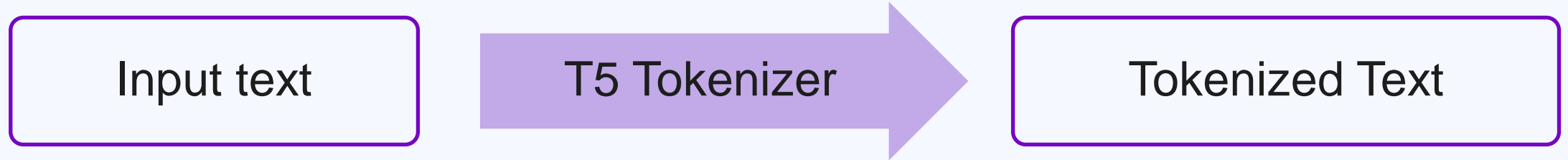
# T5: Tokenization.

- What is tokenization
  - Split text into smaller units.
  - There are different ways to do tokenization.



T5 Tokenization (Sentencepiece)

## T5: Tokenization.

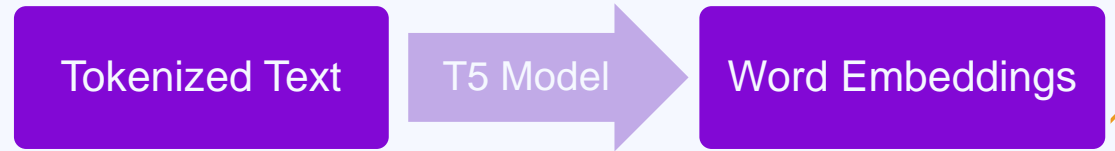


Input text: 今日はいい天気ですね。外に出て散歩したいです。

<b>_</b>	今日	は	いい	天気	ですね	。	外に出	て	散歩	したい	です	。	<b>&lt;/s&gt;</b>
<b>5</b>	4634	7	2090	18709	22034	4	29767	58	23169	5440	876	4	<b>1</b>

# Word Embeddings.

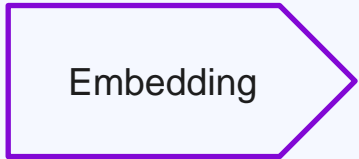
- 14 Tokens
- Each Token mapped to a vector size of 768
- 14 x 768 Vector



_	今日	は	いい	天気	ですね	。	外に出	て	散歩	したい	です	。	</s>
5	4634	7	2090	18709	22034	4	29767	58	23169	5440	876	4	1

Token Indexes

5
4634
7
...
4
1



Word Embeddings

-12.5625	12.0625	-14.0625	...	-0.2695	-5.7812	15.0625
-3.5156	-7.3750	13.6875	...	12.4375	-10.6875	0.5312
-13.3750	7.5000	-8.8125	...	3.3750	6.3125	10.1875
...	...	...	...	...	...	...
2.6250	12.1875	-10.0000	...	-0.7773	-9.7500	8.0625
-10.7500	5.6875	-12.4375	...	29.1250	4.5625	21.1250

No. of Tokens

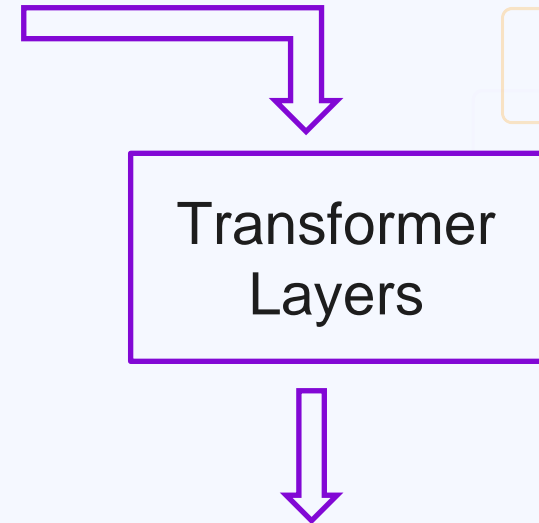
768

# Contextualized Embeddings



Word Embeddings

-12.5625	12.0625	-14.0625	...	-0.2695	-5.7812	15.0625
-3.5156	-7.3750	13.6875	...	12.4375	-10.6875	0.5312
-13.3750	7.5000	-8.8125	...	3.3750	6.3125	10.1875
...	...	...	...	...	...	...
2.6250	12.1875	-10.0000	...	-0.7773	-9.7500	8.0625
-10.7500	5.6875	-12.4375	...	29.1250	4.5625	21.1250



Contextualized Embeddings

-2.4007e-01	2.7876e-01	-3.2107e-01	...	-5.1860e-01
4.3266e-01	6.2798e-01	-8.6244e-02	...	-6.8169e-01
...	...	...	...	...
7.5833e-03	5.4663e-03	-7.6981e-03	...	-1.4981e-02

# T5 Embeddings.

- Encapsulate various aspects of the text's meaning.
- Allowing the model to perform complex text-to-text tasks.
- Form the basis of the model's ability to understand and generate text.

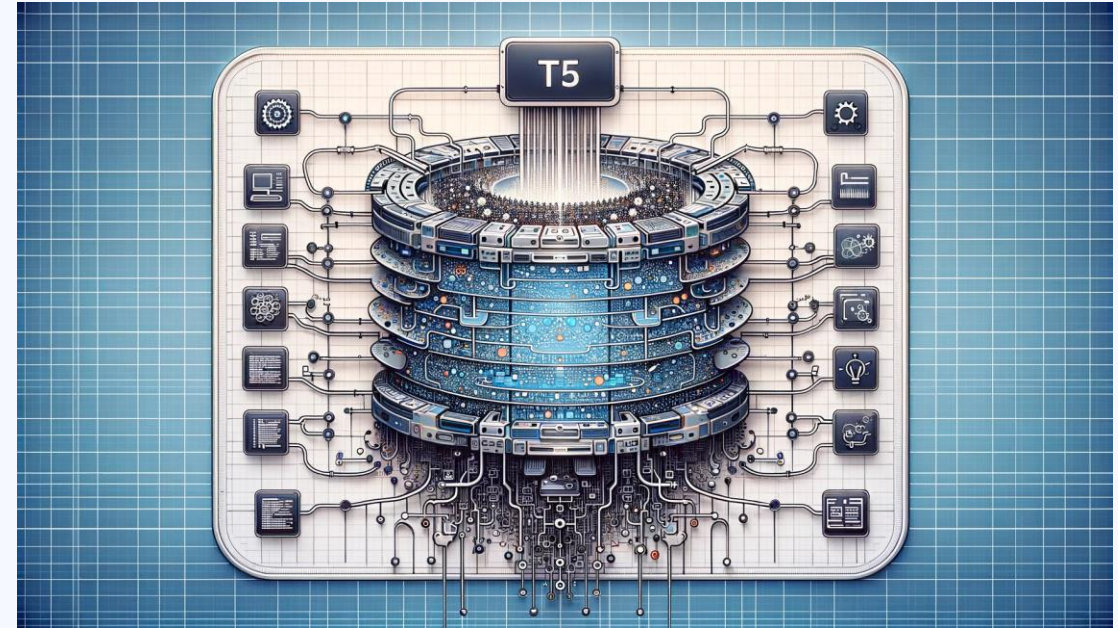
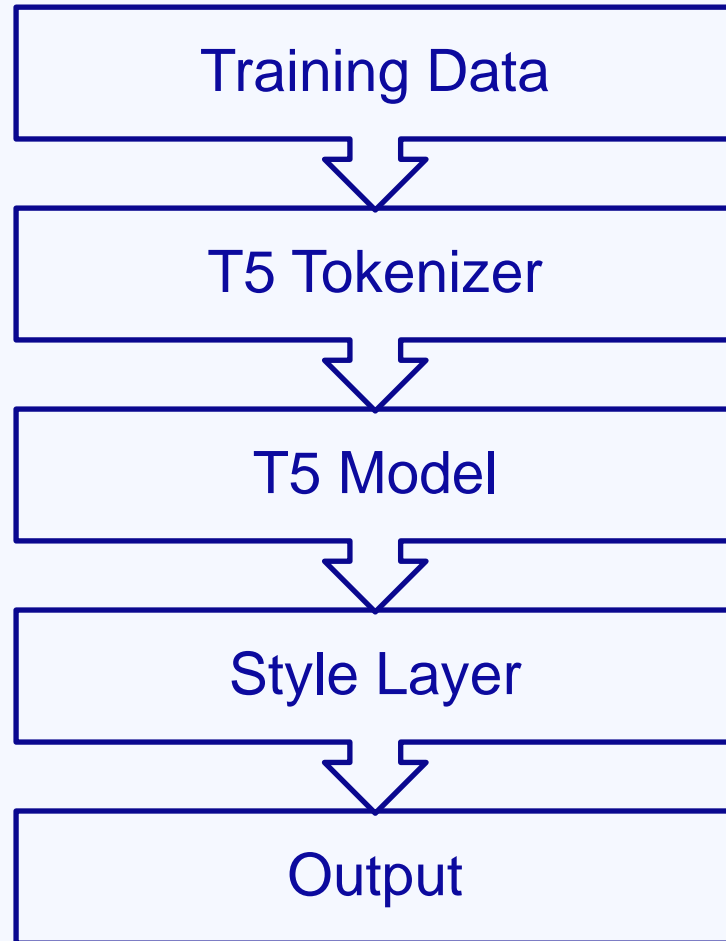
-2.4007e-01	2.7876e-01	-3.2107e-01	...	-5.1860e-01
4.3266e-01	6.2798e-01	-8.6244e-02	...	-6.8169e-01
...	...	...	...	...
7.5833e-03	5.4663e-03	-7.6981e-03	...	-1.4981e-02

# Training Data.

- Preprocess, clean.
- No need to label.
  - Subject
  - Email Body
- Emails from company employees.

	subject	email_body
0	Security Alert. Your accounts was hacked by cr...	hello!i have very bad news for you.
1	オ [redacted] 様管理画面ご説明について	クオリティア [redacted] 様お世話になっております。 [redacted] 件の件ですが、3/7の週で...
2	[redacted] について	株式会社クオリティア [redacted] 様お世話になっております。 [redacted] 以前相談させてい...
3	お見積りのご相談	クオリティア [redacted] 様お世話になっております。 [redacted] 長題の件について、現在off...
4	Active!mail アカウントその他の質問	株式会社トランスウェア [redacted] 様お世話になります。 [redacted] す。先日は評価版の送付有難う...
...	...	...
58028	Active!hunterで検知したウィルスについて	クオリティアご [redacted] 様です。いつもお世話になっております。 [redacted] ！...
58029	不審メール添付ファイル検査のお願い	クオリティアactive!hunterカスタマーサポートセンターご [redacted] 様様いつもお世話になっ...
58030	ホワイトリスト、ブラックリストの優先順位について	株式会社クオリティアサポートご [redacted] 様様いつもお世話になっております。エ [redacted] と申しま...

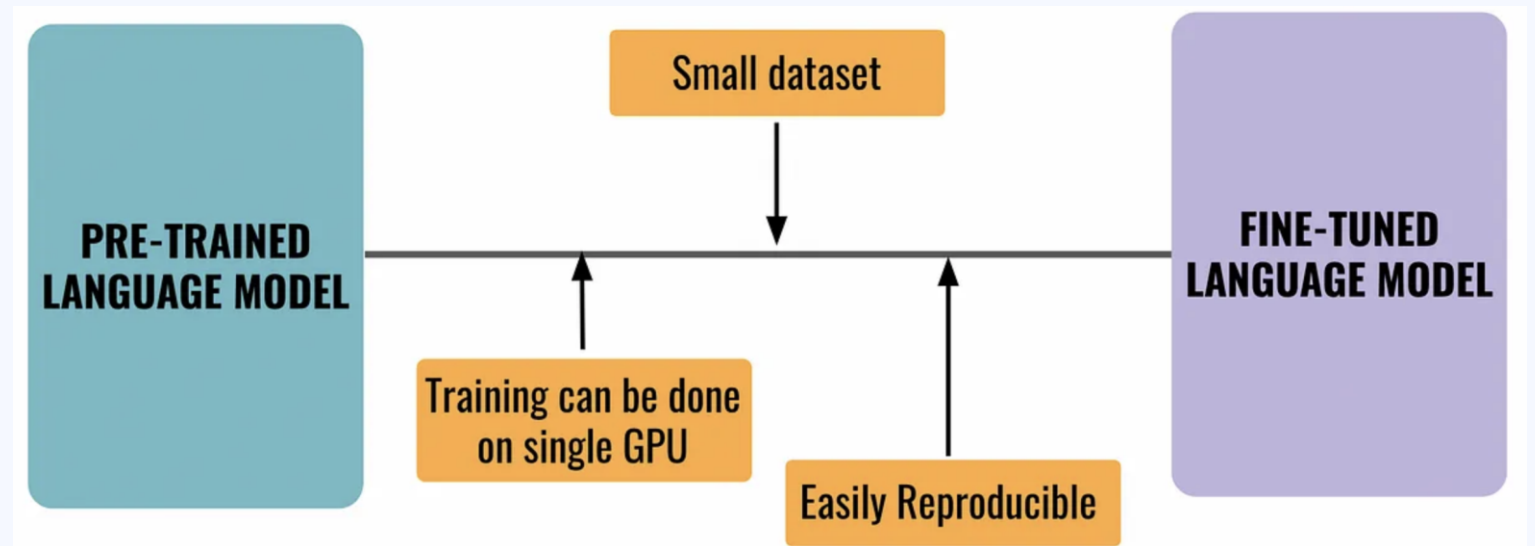
# Model Architecture.





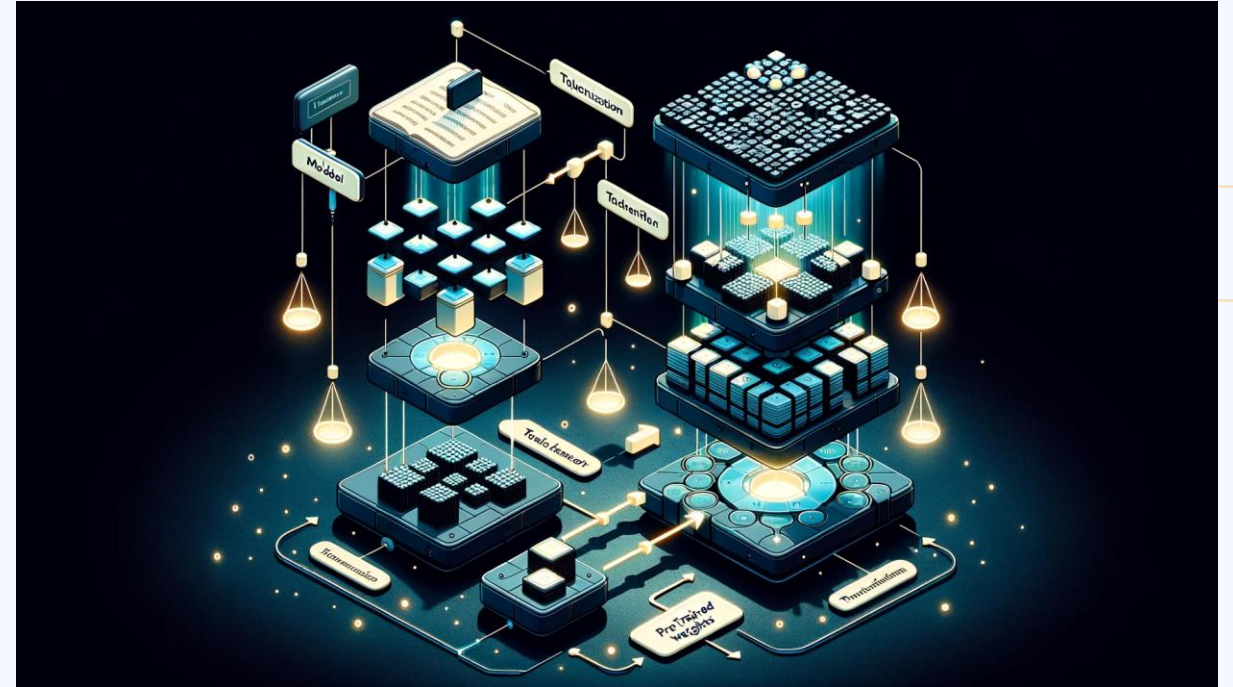
# What is Fine-tuning.

1. Process of adapting a pre-trained model to specific task.
2. Don't need to trained from scratch.
3. Efficient.
4. Better performance.
5. Flexible.



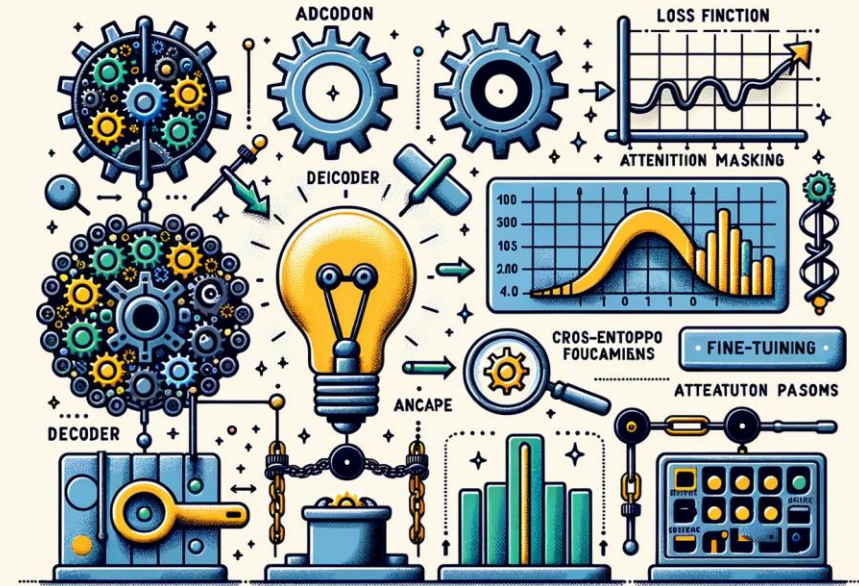
# T5: Fine-tuning Process.

1. **Tokenization**
2. **Model Architecture**
3. **Pre-trained Weights**



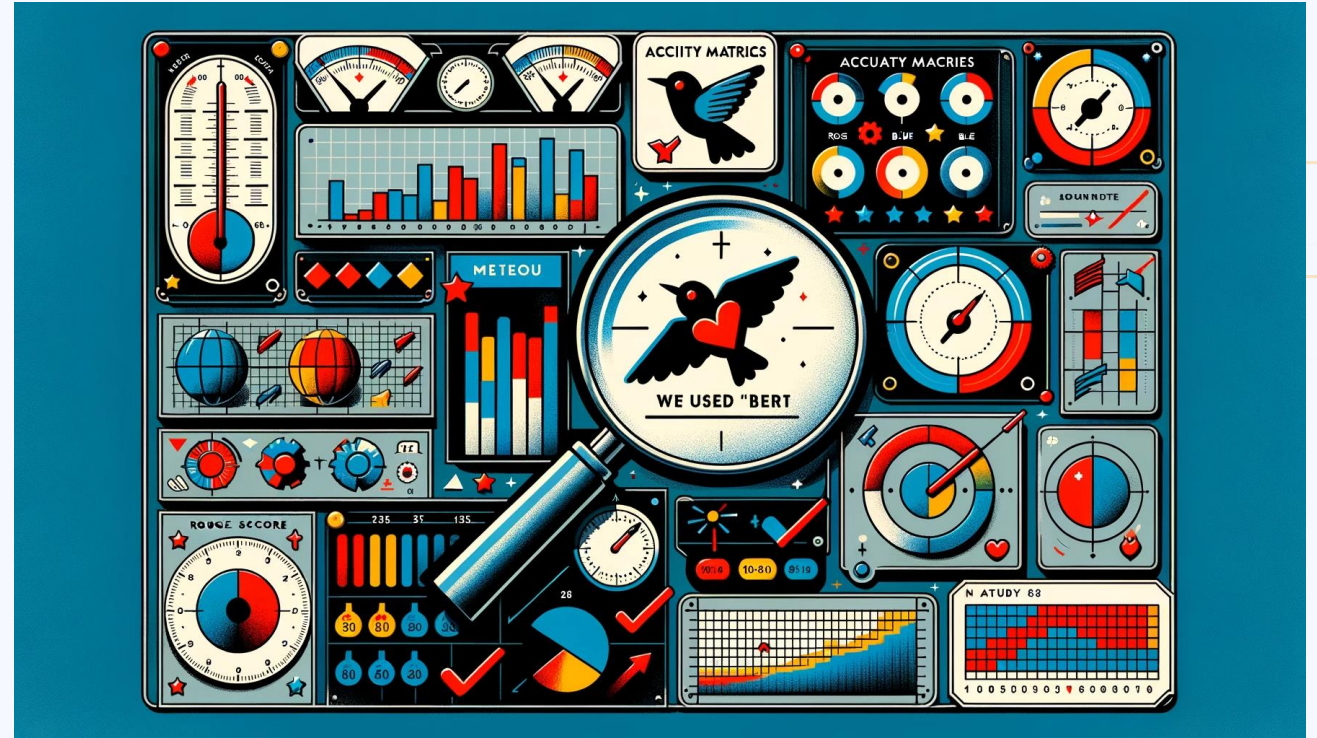
# T5: Fine-tuning Process.

4. **Decoder Adaption**
5. **Attention Masking**
6. **Loss Function**
7. **Fine-tuning Parameters**



# Accuracy Matrices.

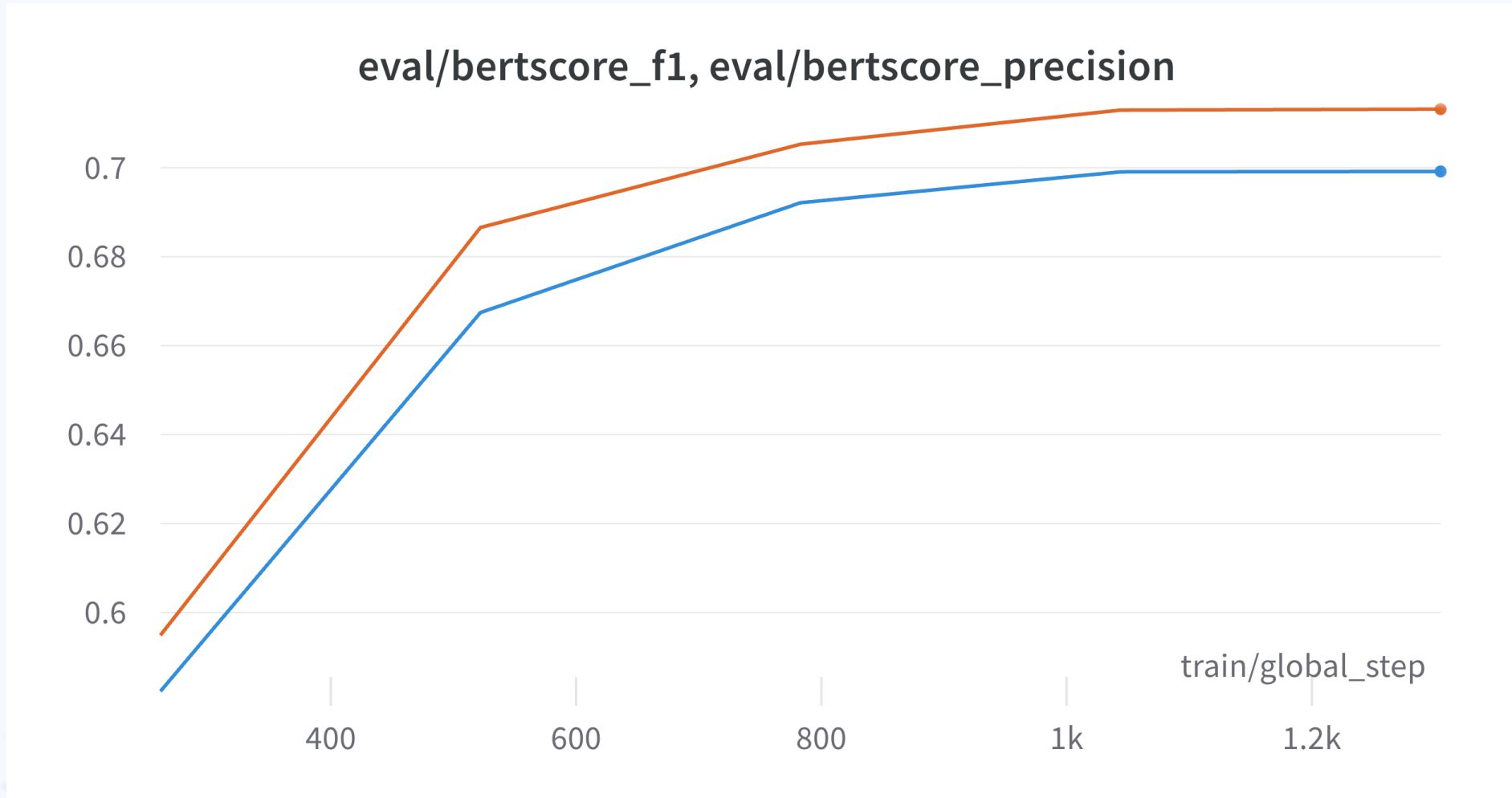
- Many Evaluation Matrices.
- Famous:
  - ROUGE Score
  - BLEU
  - METEOR
  - BERT Score
  - Etc
- We Used BERT Score.



## BERT Score.

- Uses BERT embeddings for evaluation.
- Overcomes n-gram metric limitations.
- Measures: Precision, Recall, F1-Score.
- Benefits:
  - Resilient to rephrasing.
  - Aligns with human judgments.
- Applications: Summarization, Translation.

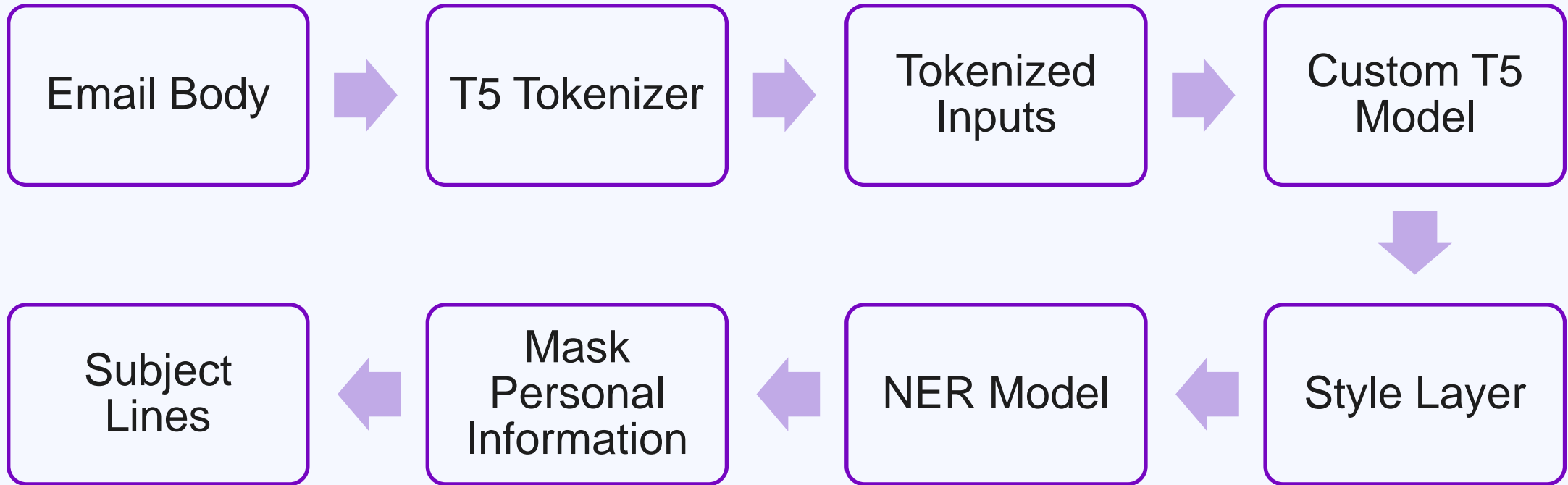
# Our Results.







# Final Pipeline.





# Subject Generation DEMO.

The screenshot displays the QUALITIA CLOUD email interface. At the top, the header includes the logo, a search bar for emails within 90 days, and user information for 'nuwan@qualitia.com'. The left sidebar shows navigation options for Email, Address book, Chat, and SSO, along with a folder list including INBOX, Sent, Trash (highlighted), Drafts (6), Spam, Reserved, Quarantine, and several archive folders. The main content area shows an empty inbox with a message: 'Emails in this folder will be automatically deleted after 90 days'. Below this, there is a 'No mail' icon consisting of a blue envelope and a document icon.

# Subject Generation Demonstration.

## □ Actual Subject: インフルエンザ予防接種 詳細



qualityia.co.jp

インフルエンザ予防接種 詳細

各位

お疲れ様です。人事総務部のです。  
標記の件、日程が決まりましたのでご連絡いたします。

期間：11月5日(火)～11月29日(金)

- ※毎週水曜日は担当医師1名のみでの営業になるため非常に混雑します。
- 水曜日の受診はなるべく避けてください。
- 上記期間内で必ず受診してください。

受付時間：

AM9: 00～12: 30

PM2: 00～4: 30

- ※毎週水曜日のみ午前中受付開始がAM9: 30～となります。
- ※上記時間内に必ず受付を完了してください。

クリニック名： クリニック

住所： 〒103-0007

東京都中央区日本橋茅場町

- 持ち物： ・添付ファイル「インフルエンザ予防接種 予診票」  
・(今回のクリニックに掛かったことのある方は)診察券  
※こちら忘れると接種できませんので事前に記入し持参してください。

費用： 全額当社負担(個別での支払いは不要です。)

その他注意点： 卵アレルギーのある方は接種できません。  
また、一般の患者さんで込み合う時間帯もあるため、  
待ち時間が長くなる場合がありますのでご注意ください。  
待合スペースは非常に狭いです。  
5名、10名単位で一緒に受診することは避けてください。

添付の「インフルエンザ予防接種 注意点」も必ずご一読をお願いします。

ご質問、ご不明点は人事総務部 までご連絡くださいませ。

皆様のご協力、何卒よろしくお願い申し上げます。

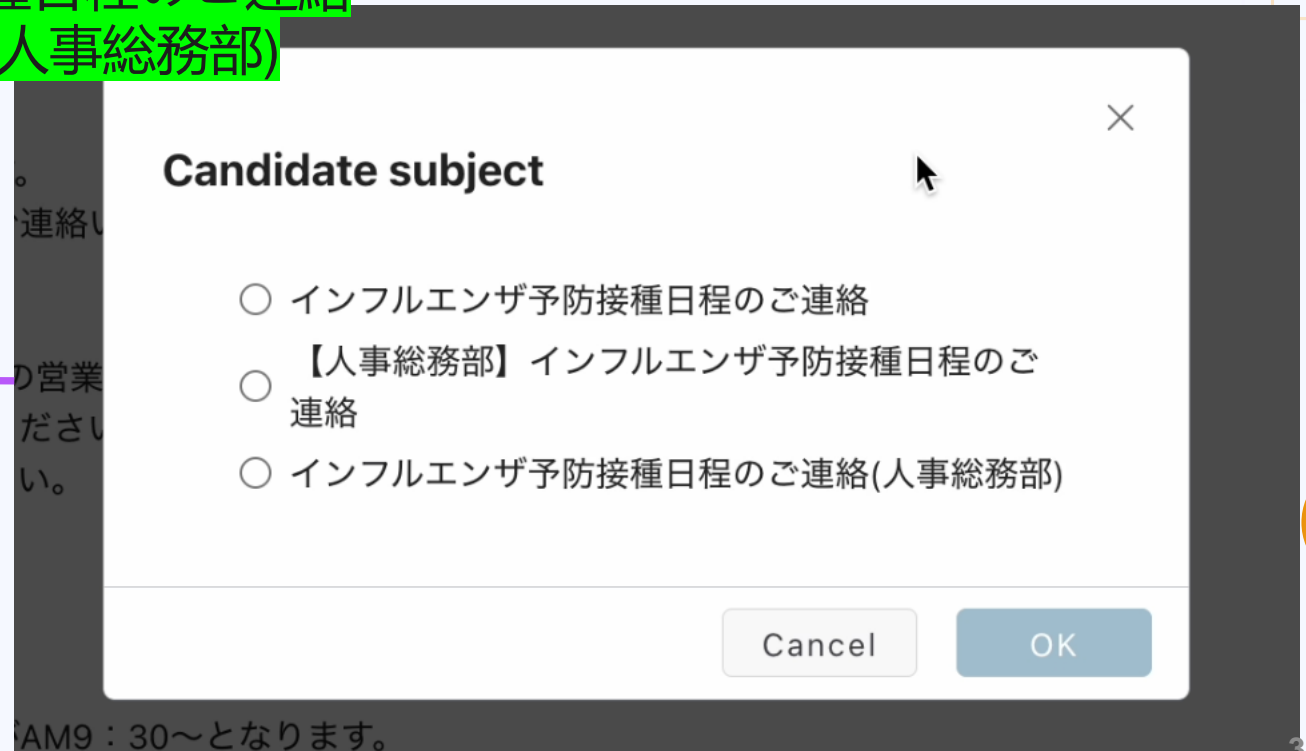
人事総務部

# Subject Generation Demonstration.

□ Actual Subject: インフルエンザ予防接種 詳細

□ Generated Subjects:

- 【インフルエンザ予防接種日程のご連絡
- 【人事総務部】インフルエンザ予防接種日程のご連絡
- インフルエンザ予防接種日程のご連絡(人事総務部)



# 04

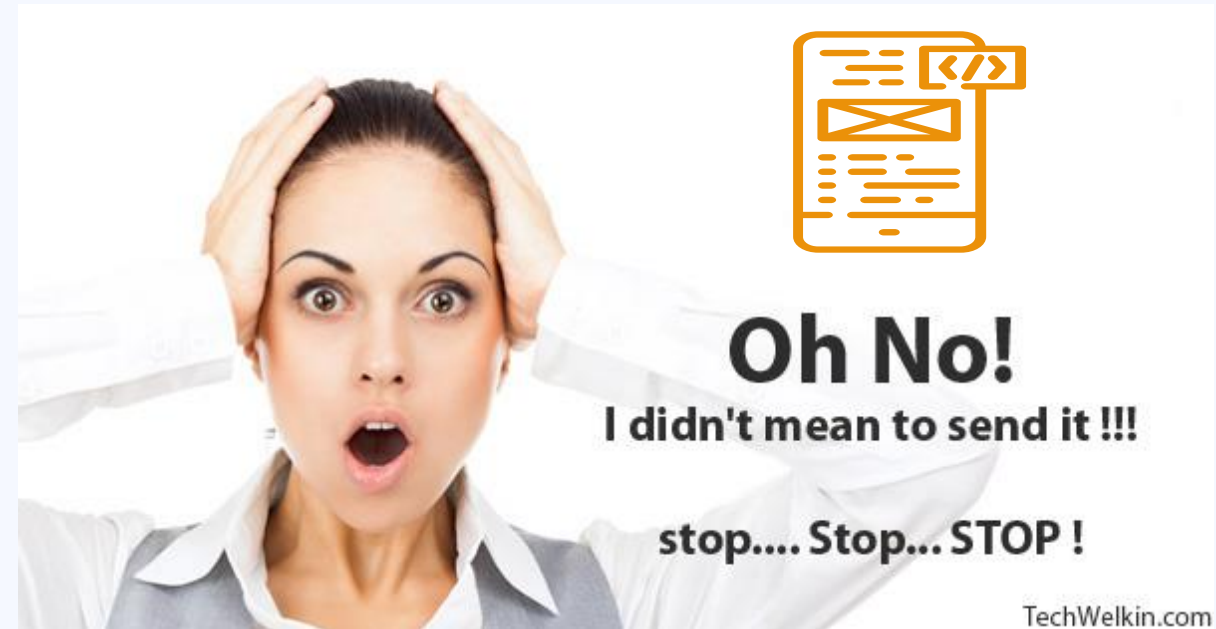
## Email Mis-sending Prevention

---

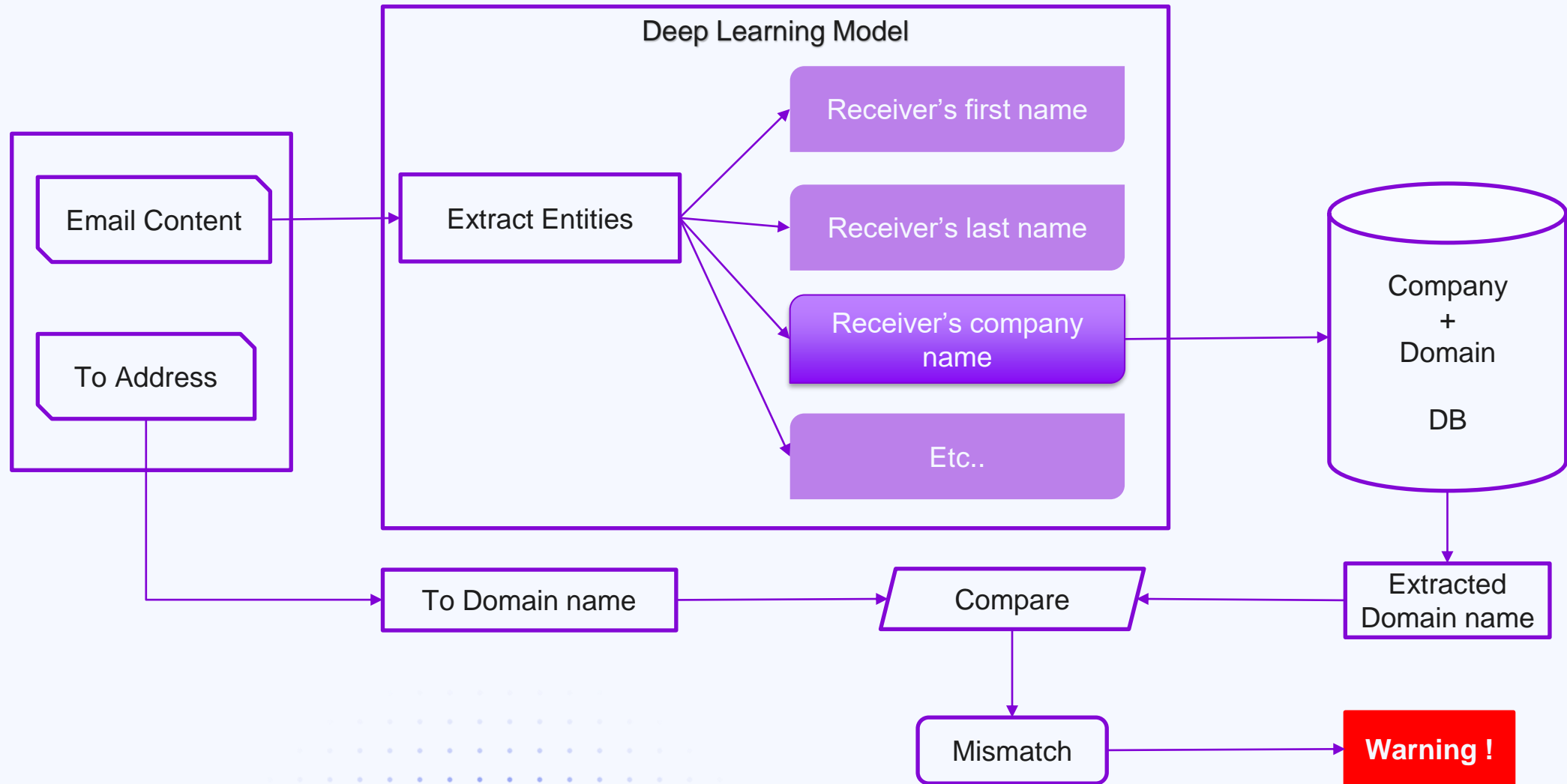
**By Named Entity Recognition.**

# Email Mis-Sending !.

- ❑ How ?
  - Incorrect email address.
  
- ❑ Damages.
  - Leaking personal information.
  - Severe consequences.
  - Leaking secrets.
  - Leaking sensitive information.
  - Damage professional relationships.



# Our Method.



# Japanese E-mails.

## Examples

シルフカンパニーのカミーユ・ビダン様、

お疲れ様です。アナハイム・エレクトロニクスのアムロ・レイと申します。

先日お送りした技術資料において、いくつかの誤りが見つかりました。誠に申し訳ございません。

修正版を準備し、明日中にはお送りさせていただきます。何卒宜しくお願い申し上げます。|

アナハイム・エレクトロニクスのアムロ・レイ様、

こんにちは、シルフカンパニーのカミーユ・ビダンと申します。来週予定されているシステムアップデートについて、技術サポートをお願いしたく思います。

アップデートの詳細とサポート可能な時間をご連絡いただければ幸いです。宜しくお願い致します。|





# Named Entity Recognition.

株式会社クオリアヌワン セネビラツナ様

# Named Entity Recognition.

株式会社クオリティアヌワン セネビラツナ様

Tokenization

# Named Entity Recognition.

株式会社クオリアヌワン セネビラツナ様

Tokenization

株式会社

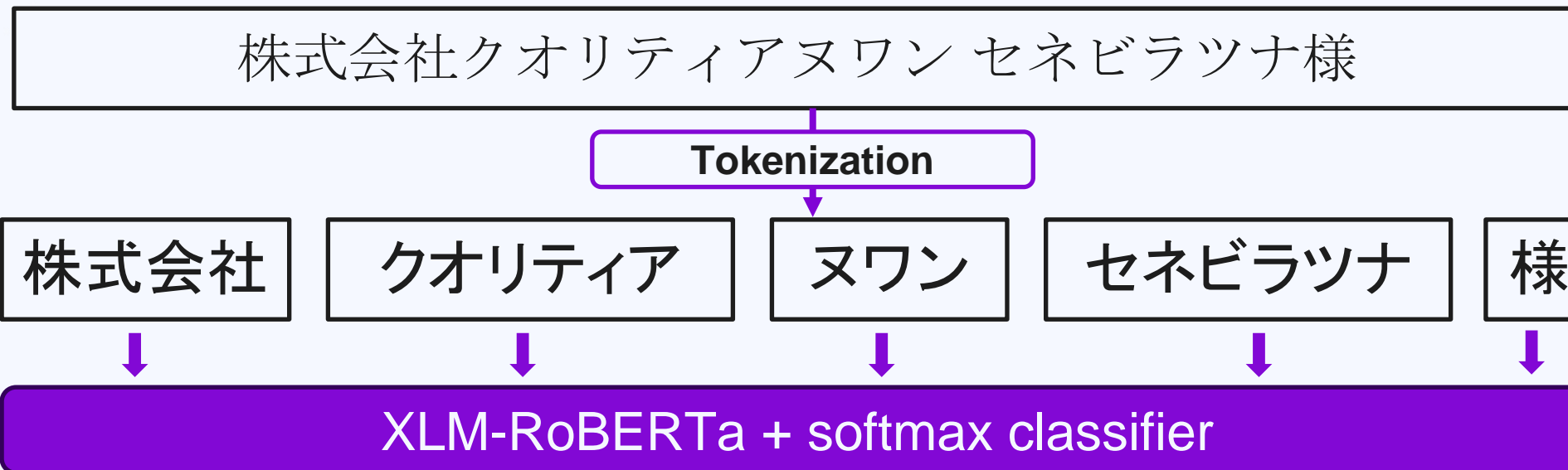
クオリア

ヌワン

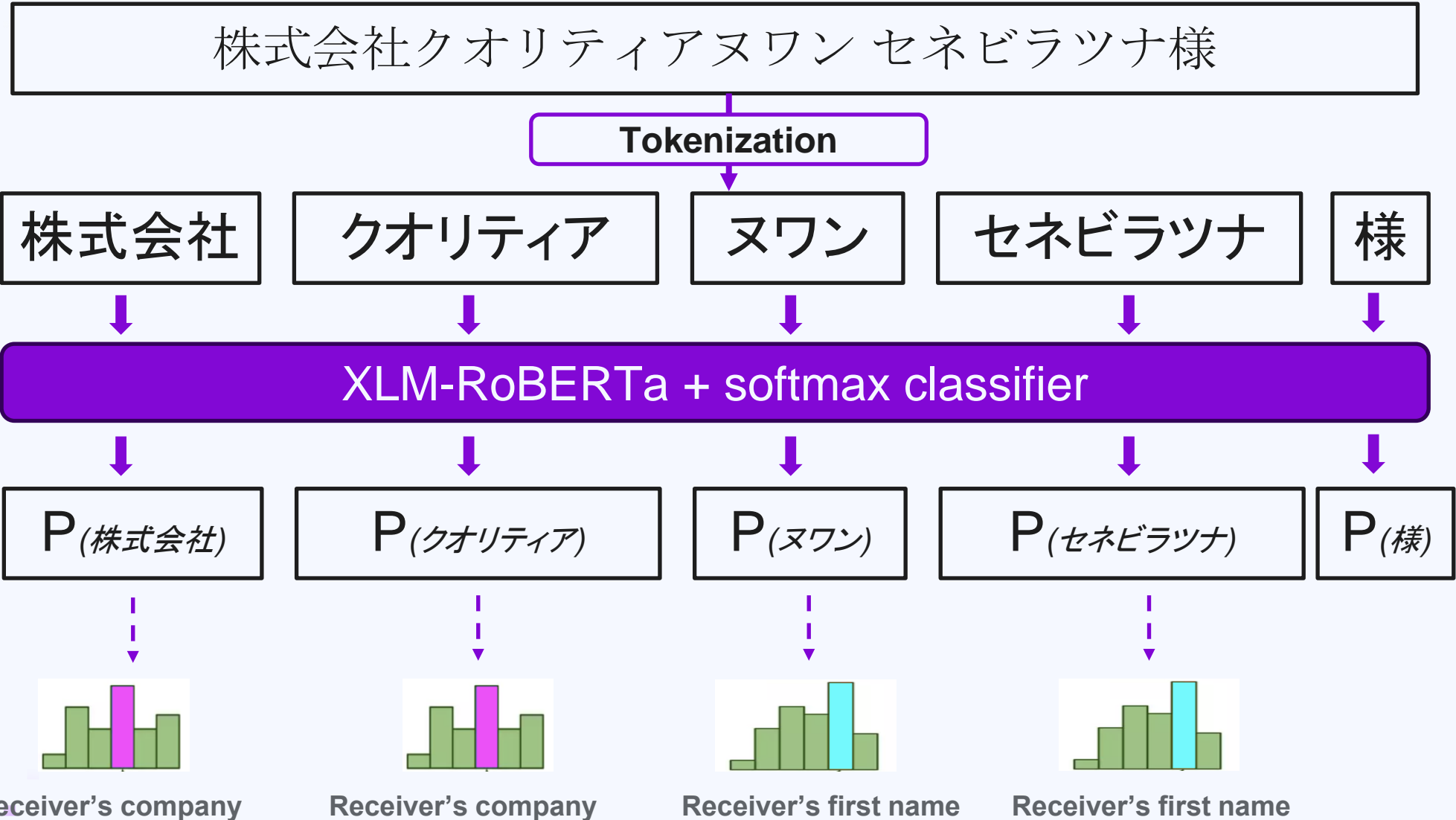
セネビラツナ

様

# Named Entity Recognition.

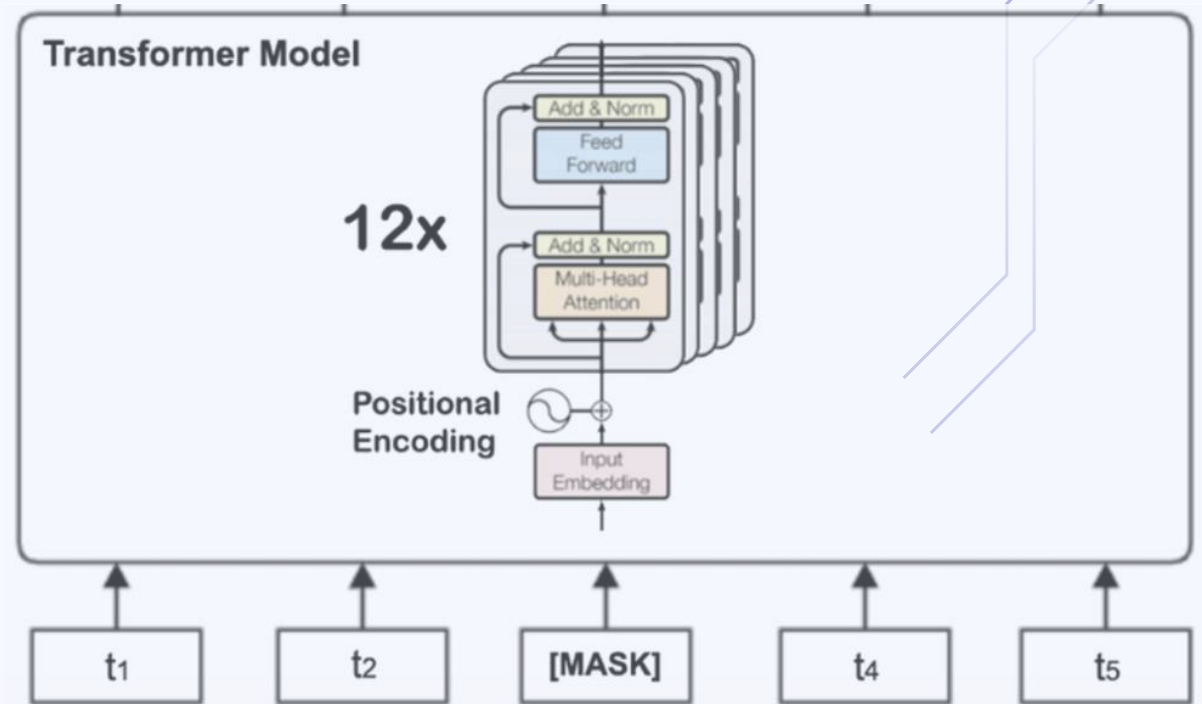


# Named Entity Recognition.



# What is XLM-RoBERTa.

- ❑ XLM-RoBERTa (Cross-lingual Model-RoBERTa).
- ❑ RoBERTa, a robustly optimized BERT variant.
- ❑ Supports 100+ languages.
- ❑ Transformer-based, like BERT.
- ❑ Encoder Only.
- ❑ 12 encoder layers for base model.

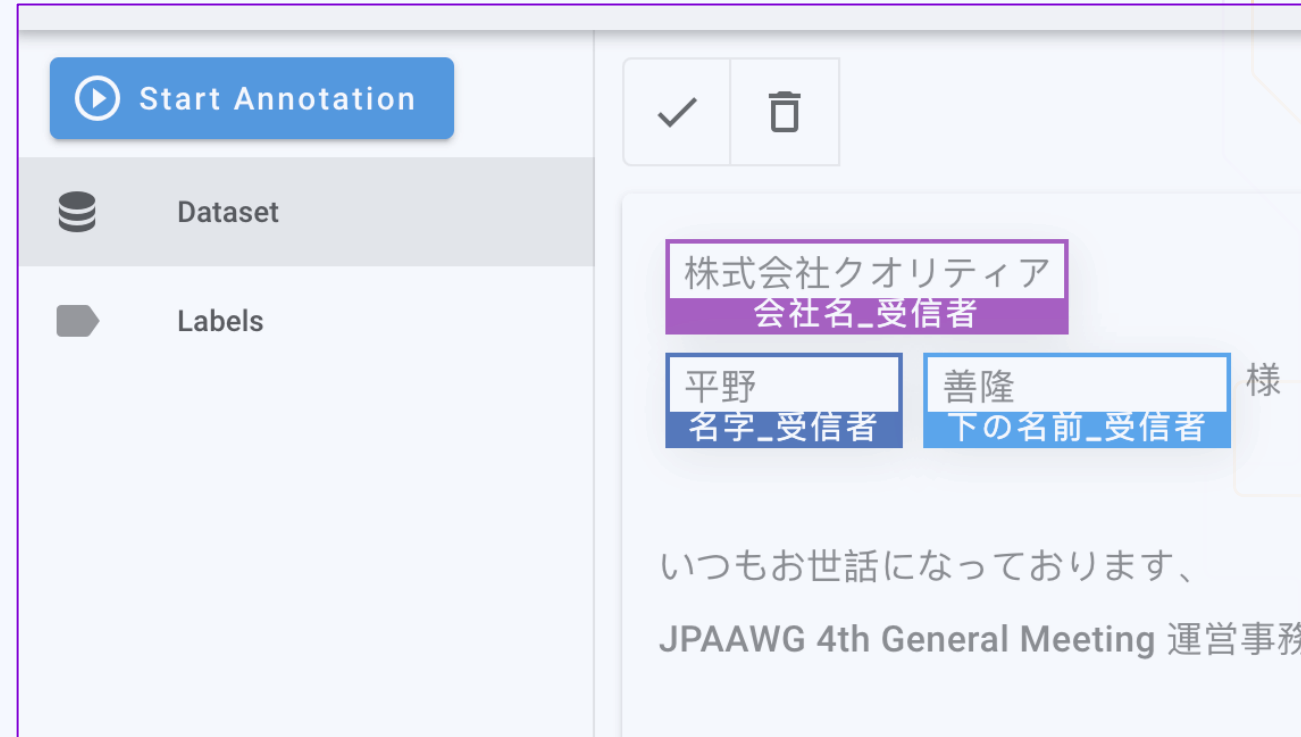




# Implement NER System.

## □ Prepare training data

- Collect Data
- Filter and Select Data
- Tag dataset.
- Doccano
  - <https://github.com/doccano/doccano>
- 26 Entities



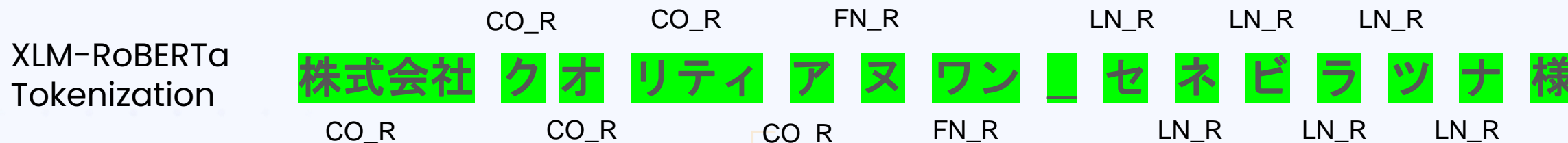
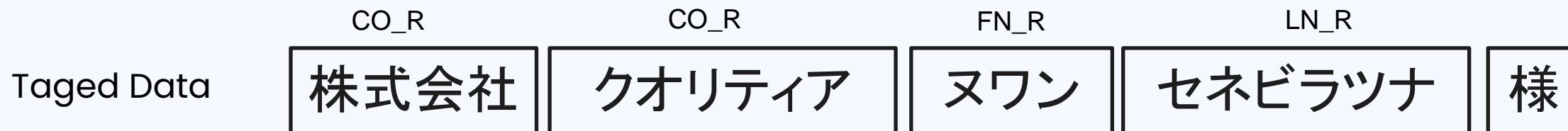
# Tagged Data.

```
1  {
2    'id': 6082,
3    'text': 'アナハイム・エレクトロニクス業ご担当者様\n\nいつもお世話になっております。 \nモルゲンレーテ社の風花と申します。 \n\n下記ご案内分につきまして\n本年',
4    'entities': [
5      {'name': 'アナハイム・エレクトロニクス業', 'span': [0, 5], 'type': '会社名_受信者'},
6      {'name': 'モルゲンレーテ社', 'span': [29, 40], 'type': '会社名_送信者'},
7      {'name': '風花', 'span': [41, 43], 'type': '名字_送信者'},
8      {'name': '幻夢コーポレーション', 'span': [115, 127], 'type': '会社名_その他'},
9      {'name': 'モルゲンレーテ社株式会社', 'span': [272, 287], 'type': '会社名_送信者'},
10     {'name': '風花', 'span': [302, 304], 'type': '名字_送信者'},
11     {'name': '月影', 'span': [304, 306], 'type': '下の名前_送信者'},
12     {'name': '〒148-0171 東京都千央区逢美台 9-18-5', 'span': [315, 353], 'type': '住所'},
13     {'name': '999-9999-9999', 'span': [358, 370], 'type': '電話番号_送信者'}
14   ]
15 }
16
```

# Fix Sequence Mismatch.

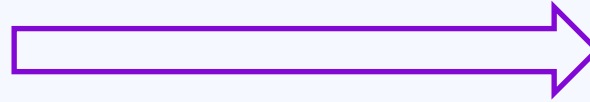
- Aligning tagged label sequence to model Tokenization.

Entity Name	Short Name
Receiver's company	CO_R
Receiver's first name	FN_R
Receiver's last name	LN_R
Other	O



# Adjusting Sequence Mismatch.

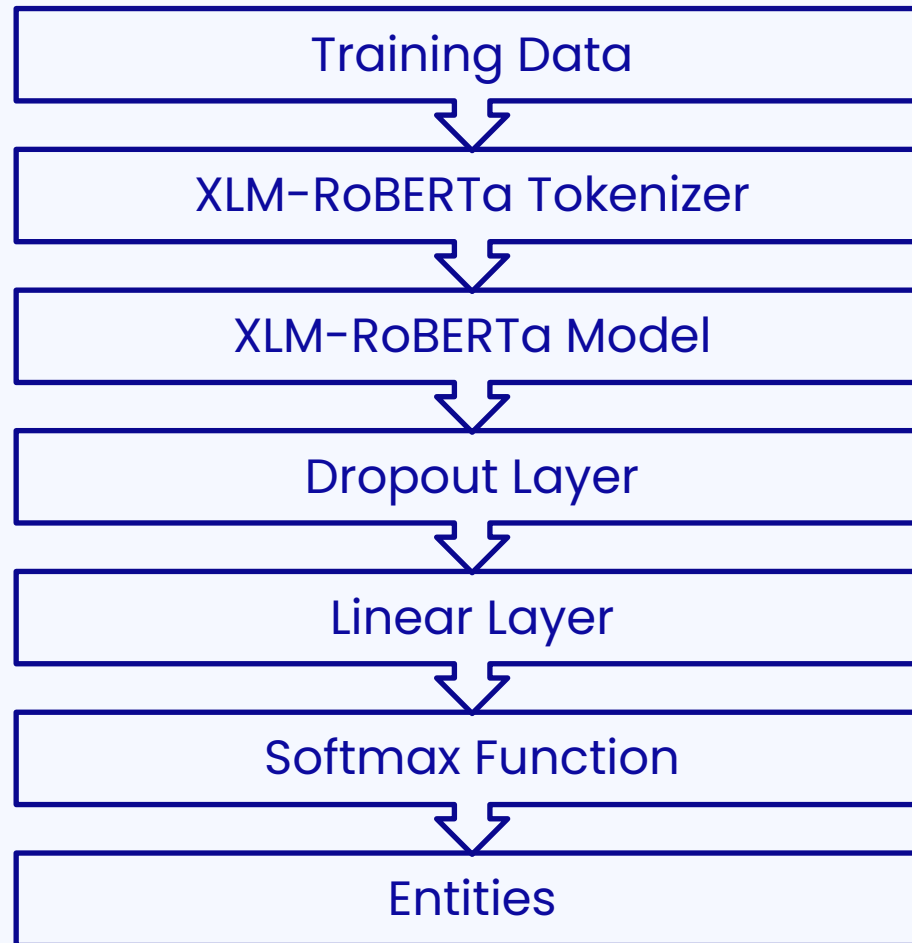
Token	Named Entity
株式会社	CO_R
クオリアティア	CO_R
ヌワン	FN_R
	O
セネビラツナ	LN_R
様	O



Token	Named Entity
株式会社	CO_R
ク	CO_R
オ	CO_R
リティ	CO_R
ア	CO_R
ヌ	FN_R
ワン	FN_R
ー	O
セ	LN_R
ネ	LN_R
ビ	LN_R
ラ	LN_R
ツ	LN_R
ナ	LN_R
様	O

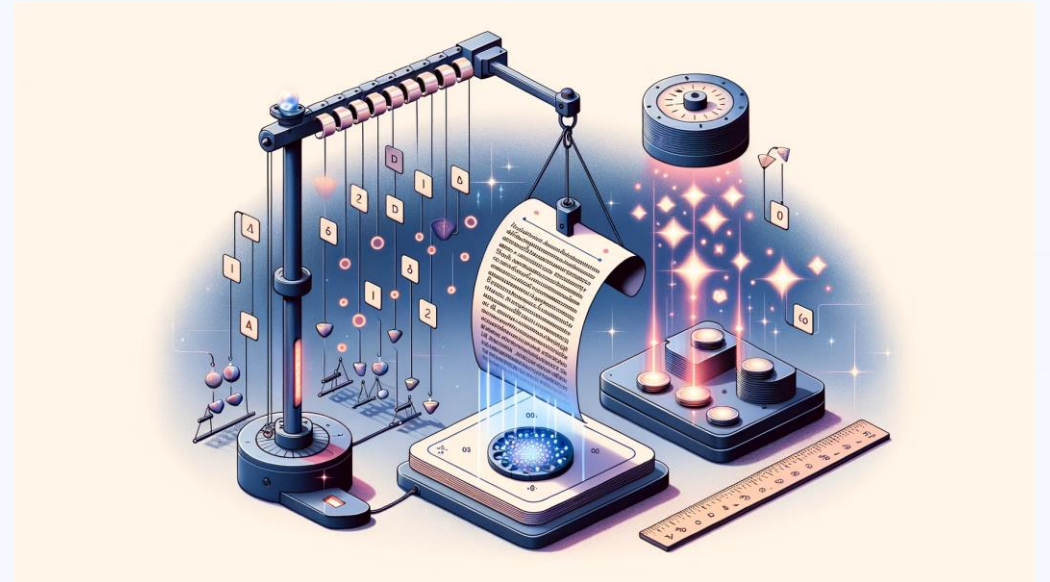
a. After adjusting for sequence mismatch for using model tokenizer.

# Model Architecture.



# Fine-tuning Process.

1. Tokenization and Text Encoding.
2. Add Special Tokens.
3. Calculate Positional Embeddings.
4. Load Pre-trained Weights.



# Fine-tuning Process.

5. Fine-tuning Layers
6. Classification Head.
7. Loss Function.
8. Fine-tuning Parameters.



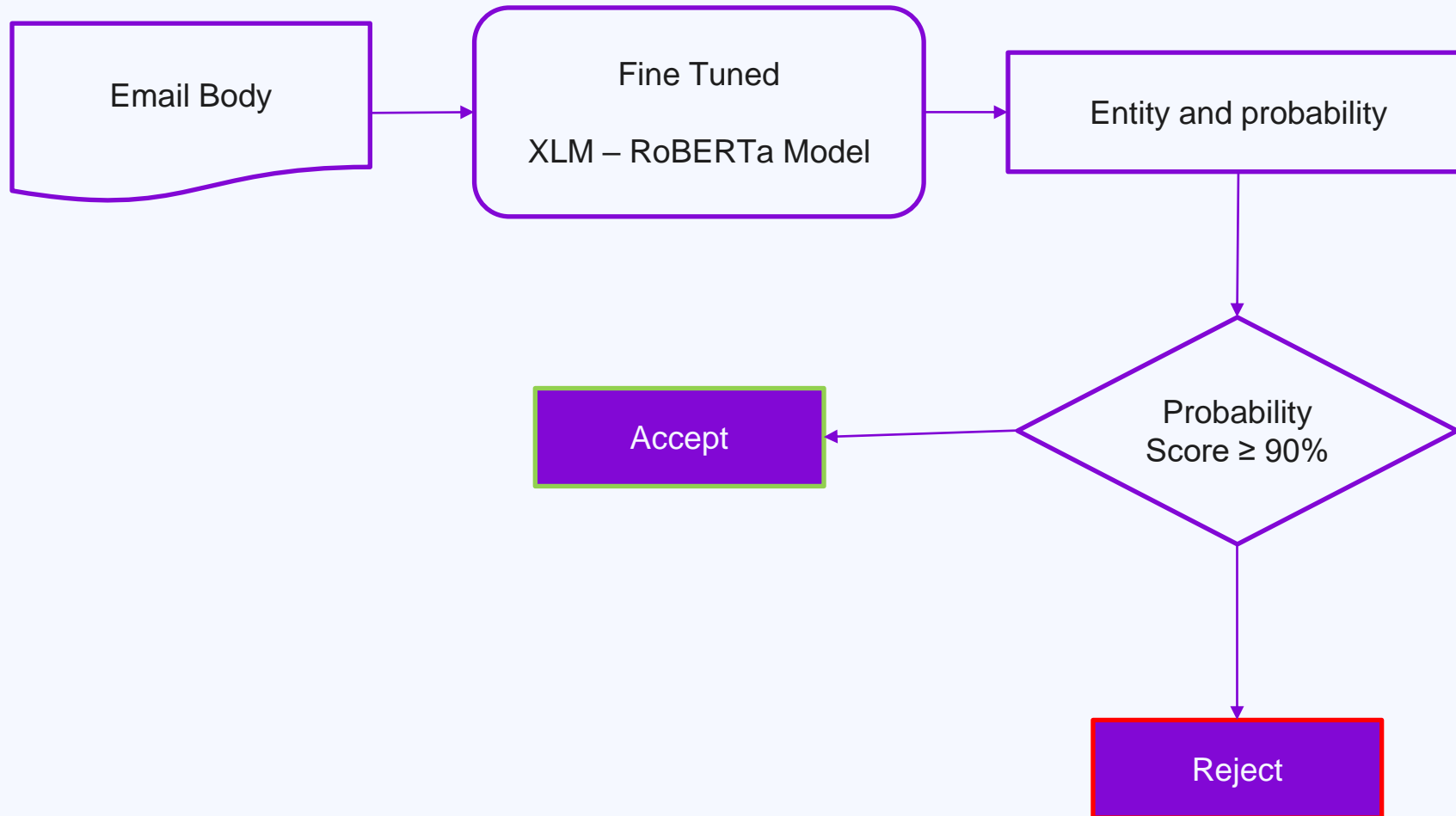


# Model Results.

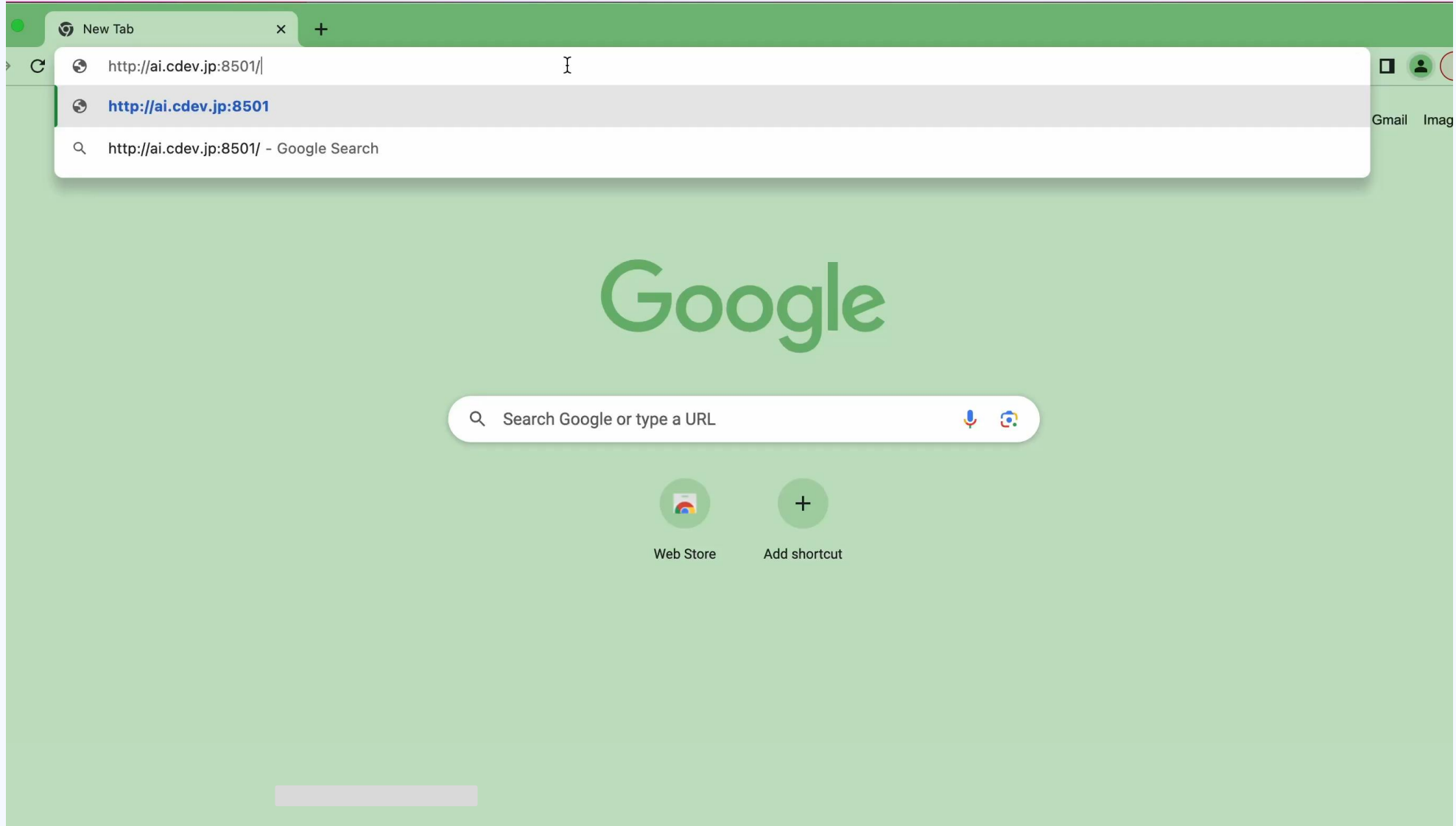
	precision	recall	f1-score
COMP-O	0.7031	0.7214	0.7121
COMP-R	0.9429	0.9445	0.9437
COMP-S	0.9229	0.9274	0.9252

- High Precision.
- High Recall.
- High F1 Score.

# Inference.




# Receiver's Company Name Extract DEMO.



# Receiver's Company Name Extract Demonstration.

## Company NER

This is a demo of a company NER models to detect the receiver's company name.

Input e-mail 

**Actual company name**

株式会社クオリティア  
平野様

いつも大変お世話になっております。  
ABCDEFGH株式会社のSenevirathneです。

=====  
※特定電子メール法第3条第1項第4号に準拠し配信しています。  
※配信停止希望は、メール配信停止フォームよりお願いします。  
<https://goo.gl/forms/XXXX>  
=====

### INFO Extracted company name

Company name: 株式会社クオリティア

Probability score: 0.9991694927215576

	Token	Tag	Token Probability
0	株式会社	COMP-R	0.9992
1	ク	COMP-R	0.9992
2	オ	COMP-R	0.9992
3	リティ	COMP-R	0.9992
4	ア	COMP-R	0.9992

# Challenges.

- Insufficient training data.
- Language and domain adaptation.
- Annotator bias and inconsistency.
- Costly to annotate data for domain specific text.

# Conclusion.

- ❖ Neural Networks.
- ❖ Transformers Architecture.
- ❖ Summarization Methods.
- ❖ T5 Model.
- ❖ Process of Building Subject Generation Model .
- ❖ Named Entity Recognition.
- ❖ XLM-RoBERTa Model.
- ❖ Building Email Mis-sending prevention Model.



# Thank You!

**Do you have any questions?**

**'The only thing that is constant is change.'**

**~ Albert Einstein ~**