# Textual Data Analysis

## Named entity recognition and sequence labeling

TURKUNLP
.ORG

UNIVERSITY
OF TURKU

# Sequence labeling

- **Input:** text, represented as sequence of tokens

- **Output:** sequence of labels (one per token) from predefined categories

  - contrast text classification: label(s) for text as a whole

- Terminology: sequence labeling can be also called *token classification* or *sequence tagging*

  - Do not confuse *sequence classification* with *token classification* or *sequence labeling*!

  - *Label*, *class* and *tag* are largely synonyms

UNIVERSITY OF TURKU

# Sequence labeling

positive

↑

**This is a good movie.**

| | |
|---|---|
| **The** | **DET** |
| **dog** | **NOUN** |
| **runs** | **VERB** |
| **in** | **ADP** |
| **the** | **DET** |
| **park** | **NOUN** |
| **.** | **PUNCT** |

**Don't copy, modify, resell, distribute or reverse engineer this app.**

↑

**Niantic grants you a limited nonexclusive nontransferable non sublicensable license to download and install a copy of the app on a mobile device and to run such copy of the app solely for your own personal noncommercial purposes. Except as expressly permitted in these terms you may not a) copy, modify or create derivative works based on the app, b) distribute transfer sublicense lease lend or rent the app to any third party, c) reverse engineer decompile or disassemble the app…**

(Example from Manor and Li, 2019.)

Text classification
Document classification
Sequence classification

Sequence labeling
Sequence tagging
Token classification

Sequence to sequence
Text generation

**TURKUNLP**
.ORG

**UNIVERSITY OF TURKU**

# Sequence labeling

- Token labels frequently represent either
    - Independent token categories (e.g. parts of speech)
    - Starts and ends of tokens spans (NER, chunking, etc.)

| The | DET |
|-----|-----|
| dog | NOUN |
| runs | VERB |
| in | PREP |
| the | DET |
| park | NOUN |
| . | PUNCT |

| New | B-GPE | } New York / GPE |
|-----|-------|---|
| York | I-GPE | |
| is | O | |
| in | O | |
| the | O | |
| United | B-GPE | } United States / GPE |
| States | I-GPE | |

TURKUNLP
.ORG

UNIVERSITY
OF TURKU

# Sequence labeling

- **Examples:**
    - Part-of-speech tagging (labels: NOUN, VERB, …)
    - Named entity recognition (labels: O, B-PER, I-PER, …)
    - Chunking (e.g. shallow parsing; labels: B-NP, I-VP, …)
    - Span marking e.g. for question answering (labels: I, O)

# Tasks: part of speech tagging

Example: Universal Dependencies POS tags

Assign each word a part of speech (POS) tag (noun, verb, etc.)

- POS tags used in older corpora varied considerably e.g. by language

- Coarse-grained "universal" tagsets common in recent work

```
ADJ: adjective
ADP: adposition
ADV: adverb
AUX: auxiliary
CCONJ: coordinating conjunction
DET: determiner
INTJ: interjection
NOUN: noun
NUM: numeral
PART: particle
PRON: pronoun
PROPN: proper noun
PUNCT: punctuation
SCONJ: subordinating conjunction
SYM: symbol
VERB: verb
X: other
```

| NOUN | NOUN | PROPN | PROPN | VERB | ADV | NOUN | PROPN | PUNCT |
|------|------|-------|-------|------|-----|------|-------|-------|
| Tasavallan | presidentti | Tarja | Halonen | matkustaa | tänään | valtiovierailulle | Irlantiin | . |

https://universaldependencies.org/u/pos/index.html

# Tasks: Named Entity Recognition

Identify **token spans** constituting **mentions of names** and assign them types

- Often extended to include also mentions of e.g. as times and dates
- Note: names frequently span **multiple tokens** (contrast POS)

Span start and extent typically marked using IOB (aka BIO) tags or a variation such as IOBES (adds [E]nd, [S]ingle)

```
Barack  B-Person
Obama   I-Person
was     O
born    O
in      O
Hawaii  B-Location
```



UNIVERSITY OF TURKU

# Named Entity Recognition: BIO tags

**Begin-In-Out** (BIO, aka IOB) tagging is frequently used to represent annotation that marks (non-overlapping) sequences of tokens

- Begin: start of annotated span
- In: token inside annotated span
- Out: not part of annotated span

Also reduced form IO (In-Out) and extended form IOBES (+End-Single)

**Example**: named entity recognition

```
Barack  B-Person
Obama   I-Person
was     O
born    O
in      O
Hawaii  B-Location
```

UNIVERSITY OF TURKU

# Tasks: Span marking / classification

IOB tagging can be applied to mark any **continuous**, **non-overlapping** spans of tokens and assign them to categories

- Phrases (chunks)
- Argumentative zones →
- Semantic roles
- Hedged claims (e.g. "may …")
- …



Figure from https://www.cl.cam.ac.uk/~sht25/az.html

# Tasks: Character sequences

Sequence labeling in NLP not limited to *token* sequences

Example: joint tokenization and sentence segmentation with labels

- **token-break**: token ends after character
- **sentence-break**: sentence ends after character
- **inside**: no break after character

<u>Note</u> that pre-trained models cannot label smaller items than their subwords!

```
INPUT: Is it you?

I       inside
s       token-break
        token-break
i       inside
t       token-break
        token-break
y       inside
o       inside
u       token-break
?       sentence-break
```

UNIVERSITY OF TURKU

# Part of Speech (POS) datasets

- Universal Dependencies: syntactic annotation for 100+ languages (https://universaldependencies.org/)



Nivre et al. (2016) Universal Dependencies v1: A Multilingual Treebank Collection

UNIVERSITY OF TURKU

# NER datasets

- CoNLL'02/03: Spanish, Dutch, English, German
  - PER(SON), LOC(ATION), ORG(ANIZATION), MISC
- OntoNotes: English, Chinese, Arabic:
  - PER, LOC, ORG, FAC(ILITY), GPE (geopolitical entity), PRODUCT, LANGUAGE, LAW, DATE, MONEY, etc.
- Turku NER corpus / TurkuONE: Finnish
  - Turku NER includes 6 entity types (PER, ORG, LOC, PRO, EVENT, DATE), TurkuONE increases data and extends to 18 types

Tjong Kim Sang and De Meulder (2003): Introduction to the CoNLL-2003 Shared Task
Hovy et al. (2006) OntoNotes: The 90% Solution
Luoma et al. (2020) A Broad-coverage Corpus for Finnish Named Entity Recognition
Luoma et al. (2021) Fine-grained Named Entity Annotation for Finnish

TURKUNLP
.ORG

UNIVERSITY
OF TURKU

# Data notebook

- https://github.com/TurkuNLP/textual-data-analysis-course/blob/main/sequence_labeling_dataset_examples.ipynb

# Sequence labeling: methods

- Rule-based systems
- (Some) Supervised machine learning approaches:
    - Hidden Markov Models
    - Conditional Random Fields
    - Convolutional and Recurrent neural nets
    - Transformers
- Zero/few-shot approaches

UNIVERSITY OF TURKU

# Recap: Classification with transformers

# Recap: Classification with transformers

- For BERT and many related models, classification head attached to special [CLS] token added to start of token sequence

    - … other options also viable, e.g. pooling (max/avg/etc.) token representations

- Additional, randomly initialized output layer added to pre-trained model

- Fine-tuning can either only train this layer (faster) or continue training also other parts of the model (higher capacity)

UNIVERSITY OF TURKU

# Recap: Classification with transformers

```python
import torch
import transformers

model_name = "bert-base-cased"
model = transformers.AutoModelForSequenceClassification.from_pretrained(model_name, num_labels=2)
```



TURKUNLP.ORG

UNIVERSITY OF TURKU

# Recap: Classification with transformers

```python
import torch
import transformers

model_name = "bert-base-cased"
model = transformers.AutoModelForSequenceClassification.from_pretrained(model_name, num_labels=2)
```

- **AutoModel** loads the correct model architecture (Bert in our case)

- **ForXXX** loads the correct output configuration ("head")

  - BertModel (no head, used e.g. for embedding)

  - ForPreTraining (loads pretraining head)

  - ForSequenceClassification

UNIVERSITY
OF TURKU

# Recap: Classification with transformers

Class logits

Classification head
(BertForSequenceClassification) →  Linear (output)

(Dropout)

Pooler output (BertModel) →  Linear + Tanh

Contextualized vector
representations of tokens
(BertModel) →

| $T_0$ | $T_1$ | $T_2$ | $T_3$ | $T_4$ | $T_5$ | ... | $T_N$ |

Transformer

| $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | ... | $E_N$ |

| [CLS] | This | is | a | good | movie | ... | |

TURKUNLP
.ORG

UNIVERSITY
OF TURKU

# Recap: Classification with transformers

```python
class BertForSequenceClassification(BertPreTrainedModel):
    def __init__(self, config):
        super().__init__(config)
        self.num_labels = config.num_labels
        self.config = config

        self.bert = BertModel(config)
        classifier_dropout = (
            config.classifier_dropout if config.classifier_dropout is not None else config.hidden_dropout_prob
        )
        self.dropout = nn.Dropout(classifier_dropout)
        self.classifier = nn.Linear(config.hidden_size, config.num_labels)

        # Initialize weights and apply final processing
        self.post_init()
```

Bert model without any output layer

Optional dropout

New linear output layer

TURKUNLP
.ORG

UNIVERSITY
OF TURKU

```python
def forward(
    self,
    input_ids: Optional[torch.Tensor] = None,
    attention_mask: Optional[torch.Tensor] = None,
    token_type_ids: Optional[torch.Tensor] = None,
    position_ids: Optional[torch.Tensor] = None,
    head_mask: Optional[torch.Tensor] = None,
    inputs_embeds: Optional[torch.Tensor] = None,
    labels: Optional[torch.Tensor] = None,
    output_attentions: Optional[bool] = None,
    output_hidden_states: Optional[bool] = None,
    return_dict: Optional[bool] = None,
) -> Union[Tuple[torch.Tensor], SequenceClassifierOutput]:
    r"""
    labels (`torch.LongTensor` of shape `(batch_size,)`, *optional*):
        Labels for computing the sequence classification/regression loss. Indices should be in `[0, ...,
        config.num_labels - 1]`. If `config.num_labels == 1` a regression loss is computed (Mean-Square
        `config.num_labels > 1` a classification loss is computed (Cross-Entropy).
    """
    return_dict = return_dict if return_dict is not None else self.config.use_return_dict

    outputs = self.bert(
        input_ids,
        attention_mask=attention_mask,
        token_type_ids=token_type_ids,
        position_ids=position_ids,
        head_mask=head_mask,
        inputs_embeds=inputs_embeds,
        output_attentions=output_attentions,
        output_hidden_states=output_hidden_states,
        return_dict=return_dict,
    )

    pooled_output = outputs[1]

    pooled_output = self.dropout(pooled_output)
    logits = self.classifier(pooled_output)
```

Bert outputs (contextualized vectors)

Extract "pooled output" from outputs

Run through dropout, and output layer to get classification logits

TURKUNLP
.ORG

```python
def forward(
    self,
    input_ids: Optional[torch.Tensor] = None,
    attention_mask: Optional[torch.Tensor] = None,
    token_type_ids: Optional[torch.Tensor] = None,
    position_ids: Optional[torch.Tensor] = None,
    head_mask: Optional[torch.Tensor] = None,
    inputs_embeds: Optional[torch.Tensor] = None,
    labels: Optional[torch.Tensor] = None,
    output_attentions: Optional[bool] = None,
    output_hidden_states: Optional[bool] = None,
    return_dict: Optional[bool] = None,
) -> Union[Tuple[torch.Tensor], SequenceClassifierOutput]:
    r"""
    labels (`torch.LongTensor` of shape `(batch_size,)`, *optional*):
        Labels for computing the sequence classification/regression loss. Indices should be in `[0, ...,
        config.num_labels - 1]`. If `config.num_labels == 1` a regression loss is computed (Mean-Square
        `config.num_labels > 1` a classification loss is computed (Cross-Entropy).
    """
    return_dict = return_dict if return_dict is not None else self.config.use_return_dict

    outputs = self.bert(
        input_ids,
        attention_mask=attention_mask,
        token_type_ids=token_type_ids,
        position_ids=position_ids,
        head_mask=head_mask,
        inputs_embeds=inputs_embeds,
        output_attentions=output_attentions,
        output_hidden_states=output_hidden_states,
        return_dict=return_dict,
    )

    pooled_output = outputs[1]

    pooled_output = self.dropout(pooled_output)
    logits = self.classifier(pooled_output)
```

Bert outputs (contextualized vectors)

What is "pooled output" (outputs

Extract "pooled output" from outputs

Run through dropout, and output
layer to get classification logits

TURKUNLP
.ORG

# Recap: Classification with transformers

class transformers.**BertModel**

( config, add_pooling_layer = T

**Parameters**

- **config** (BertConfig) — Model configu
  config file does not load the weight
  from_pretrained() method to load t

The bare Bert Model transformer outp

**Returns** transformers.modeling_outputs.BaseModelOutputWithPoolingAndCrossAttentions **or** tuple(torch.FloatTensor)

A transformers.modeling_outputs.BaseModelOutputWithPoolingAndCrossAttentions or a tuple of torch.FloatTensor (if return_dict=False is passed or when config.return_dict=False) comprising various elements depending on the configuration (BertConfig) and inputs.

- **last_hidden_state** (torch.FloatTensor of shape (batch_size, sequence_length, hidden_size)) — Sequence of hidden-states at the output of the last layer of the model.

- **pooler_output** (torch.FloatTensor of shape (batch_size, hidden_size)) — Last layer hidden-state of the first token of the sequence (classification token) after further processing through the layers used for the auxiliary pretraining task. E.g. for BERT-family of models, this returns the classification token after processing through a linear layer and a tanh activation function. The linear layer weights are trained from the next sentence prediction (classification) objective during pretraining.

- **hidden_states** (tuple(torch.FloatTensor), *optional*, returned when output_hidden_states=True is passed or when config.output_hidden_states=True) — Tuple of torch.FloatTensor (one for the output of the embeddings, if the model has an embedding layer, + one for the output of each layer) of shape (batch_size, sequence_length, hidden_size).

  Hidden-states of the model at the output of each layer plus the optional initial embedding outputs.

TURKUNLP.ORG

UNIVERSITY OF TURKU

# Recap: Classification with transformers

class transformers.**BertModel**

( config, add_pooling_layer = Tr

**Parameters**

- **config** (BertConfig) — Model configu
  config file does not load the weight
  from_pretrained() method to load t

The bare Bert Model transformer outp

**Returns**

transformers.modeling_outputs.BaseModelOutputWithPoolingAndCrossAttentions or
tuple(torch.FloatTensor)

A transformers.modeling_outputs.BaseM
torch.FloatTensor (if return_dict=F
various elements depending on the confi

- **last_hidden_state** (torch.FloatTe
  — Sequence of hidden-states at the

- **pooler_output** (torch.FloatTens
  state of the first token of the sequen
  used for the auxiliary pretraining tas
  after processing through a linear lay
  from the next sentence prediction (c

- **hidden_states** (tuple(torch.FloatTensor), optional, returned when output_hidden_states=True
  is passed or when config.output_hidden_states=True) — Tuple of torch.FloatTensor (one for the
  output of the embeddings, if the model has an embedding layer, + one for the output of each layer) of
  shape (batch_size, sequence_length, hidden_size).

  Hidden-states of the model at the output of each layer plus the optional initial embedding outputs.

```python
class BertPooler(nn.Module):
    def __init__(self, config):
        super().__init__()
        self.dense = nn.Linear(config.hidden_size, config.hidden_size)
        self.activation = nn.Tanh()

    def forward(self, hidden_states: torch.Tensor) -> torch.Tensor:
        # We "pool" the model by simply taking the hidden state corresponding
        # to the first token.
        first_token_tensor = hidden_states[:, 0]
        pooled_output = self.dense(first_token_tensor)
        pooled_output = self.activation(pooled_output)
        return pooled_output
```

TURKUNLP
.ORG

UNIVERSITY
OF TURKU

# Sequence labeling with transformers

| | |
|---|---|
| The | DET |
| dog | NOUN |
| runs | VERB |
| in | PREP |
| the | DET |
| park | NOUN |
| . | PUNCT |

| | | |
|---|---|---|
| New | B-GPE | } New York / GPE |
| York | I-GPE | |
| is | O | |
| in | O | |
| the | O | |
| United | B-GPE | } United States / GPE |
| States | I-GPE | |

# Sequence labeling with transformers

# Sequence labeling with transformers

# Sequence labeling with transformers

- Randomly initialized output layer attached to each token (shared weights).

- (Possible special tokens such as [CLS] generally not used)

- As in classification, fine-tuning can be performed either only for output layer or for whole model



TURKUNLP
.ORG

# Sequence labeling with transformers

```python
import torch
import transformers

model = 'bert-base-cased'
model = transformers.AutoModelForTokenClassification.from_pretrained(model, num_labels=6)
```

UNIVERSITY OF TURKU

# Sequence labeling with transformers

Classification head
(BertForTokenClassification)

Contextualized vector
representations of tokens
(BertModel)

Same layer
repeated / shared
weights!

Logits    Logits    Logits    Logits    Logits    Logits

| Linear | Linear | Linear | Linear | Linear | Linear |

| (Dropout) | (Dropout) | (Dropout) | (Dropout) | (Dropout) | (Dropout) |

$T_0$    $T_1$    $T_2$    $T_3$    $T_4$    $T_5$    ...    $T_N$

Transformer block

Transformer block

Transformer block

$E_0$    $E_1$    $E_2$    $E_3$    $E_4$    $E_5$    ...    $E_N$

[CLS]    I    own    a    Nokia    phone    ...

TURKUNLP
.ORG

UNIVERSITY
OF TURKU

# What happens in the code?

```python
class BertForTokenClassification(BertPreTrainedModel):
    def __init__(self, config):
        super().__init__(config)
        self.num_labels = config.num_labels

        self.bert = BertModel(config, add_pooling_layer=False)
        classifier_dropout = (
            config.classifier_dropout if config.classifier_dropout is not None else config.hidden_dropout
        )
        self.dropout = nn.Dropout(classifier_dropout)
        self.classifier = nn.Linear(config.hidden_size, config.num_labels)

        # Initialize weights and apply final processing
        self.post_init()
```

Bert model without any output layer

Optional dropout

New linear output layer

Exactly same init() as in sequence classification, except pooling layer not needed!

```python
def forward(
    self,
    input_ids: Optional[torch.Tensor] = None,
    attention_mask: Optional[torch.Tensor] = None,
    token_type_ids: Optional[torch.Tensor] = None,
    position_ids: Optional[torch.Tensor] = None,
    head_mask: Optional[torch.Tensor] = None,
    inputs_embeds: Optional[torch.Tensor] = None,
    labels: Optional[torch.Tensor] = None,
    output_attentions: Optional[bool] = None,
    output_hidden_states: Optional[bool] = None,
    return_dict: Optional[bool] = None,
) -> Union[Tuple[torch.Tensor], TokenClassifierOutput]:
    r"""
    labels (`torch.LongTensor` of shape `(batch_size, sequence_length)`, *optional*):
        Labels for computing the token classification loss. Indices should be in `[0, ..., config.num_labels - 1]`.
    """
    return_dict = return_dict if return_dict is not None else self.config.use_return_dict

    outputs = self.bert(
        input_ids,
        attention_mask=attention_mask,
        token_type_ids=token_type_ids,
        position_ids=position_ids,
        head_mask=head_mask,
        inputs_embeds=inputs_embeds,
        output_attentions=output_attentions,
        output_hidden_states=output_hidden_states,
        return_dict=return_dict,
    )

    sequence_output = outputs[0]

    sequence_output = self.dropout(sequence_output)
    logits = self.classifier(sequence_output)
```

Bert outputs (contextualized vectors)

Extract "sequence outputs"

Run through dropout, and output layer to get classification logits

TURKUNLP
.ORG

# What happens in code?

class transformers.**BertModel**

( config, add_pooling_layer = Tr...

**Parameters**

- **config** ([BertConfig](#)) — Model config...
  config file does not load the weight...
  [from_pretrained](#)() method to load t...

The bare Bert Model transformer outp...

**Returns**  [transformers.modeling_outputs.BaseModelOutputWithPoolingAndCrossAttentions](#) **or** `tuple(torch.FloatTensor)`
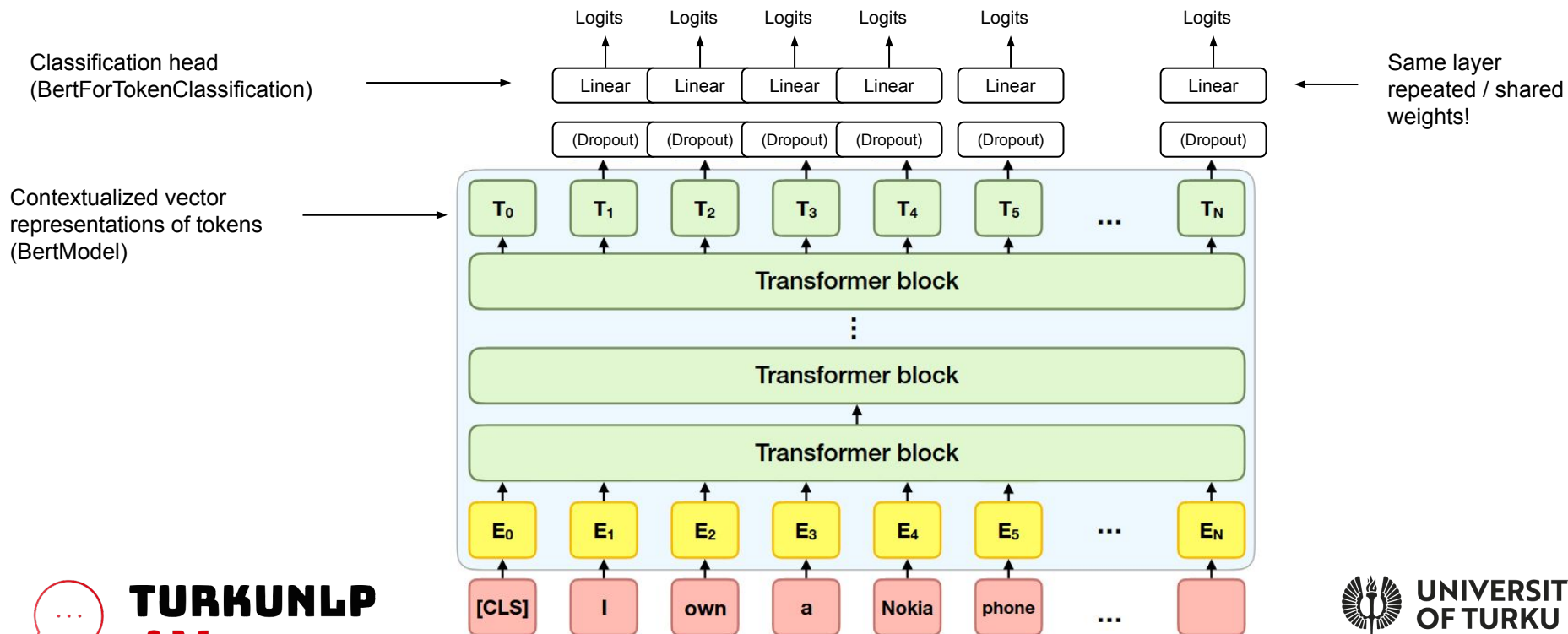
A [transformers.modeling_outputs.BaseModelOutputWithPoolingAndCrossAttentions](#) or a tuple of `torch.FloatTensor` (if `return_dict=False` is passed or when `config.return_dict=False`) comprising various elements depending on the configuration ([BertConfig](#)) and inputs.

- **last_hidden_state** (`torch.FloatTensor` of shape `(batch_size, sequence_length, hidden_size)`) — Sequence of hidden-states at the output of the last layer of the model.

- **pooler_output** (`torch.FloatTensor` of shape `(batch_size, hidden_size)`) — Last layer hidden-state of the first token of the sequence (classification token) after further processing through the layers used for the auxiliary pretraining task. E.g. for BERT-family of models, this returns the classification token after processing through a linear layer and a tanh activation function. The linear layer weights are trained from the next sentence prediction (classification) objective during pretraining.

- **hidden_states** (`tuple(torch.FloatTensor)`, *optional*, returned when `output_hidden_states=True` is passed or when `config.output_hidden_states=True`) — Tuple of `torch.FloatTensor` (one for the output of the embeddings, if the model has an embedding layer, + one for the output of each layer) of shape `(batch_size, sequence_length, hidden_size)`.

  Hidden-states of the model at the output of each layer plus the optional initial embedding outputs.

TURKUNLP .ORG

UNIVERSITY OF TURKU

# Sequence labeling with transformers

Classification head
(BertForTokenClassification)

Contextualized vector
representations of tokens
(BertModel)

Same layer
repeated / shared
weights!

Logits · Logits · Logits · Logits · Logits · Logits

Linear · Linear · Linear · Linear · Linear · Linear

(Dropout) · (Dropout) · (Dropout) · (Dropout) · (Dropout) · (Dropout)

$T_0$ · $T_1$ · $T_2$ · $T_3$ · $T_4$ · $T_5$ · ... · $T_N$

Transformer block

Transformer block

Transformer block

$E_0$ · $E_1$ · $E_2$ · $E_3$ · $E_4$ · $E_5$ · ... · $E_N$

[CLS] · I · own · a · Nokia · phone · ...

TURKUNLP
.ORG

UNIVERSITY
OF TURKU

# Sequence labeling with transformers

- Note that label sequences are not directly modeled in AutoModelForTokenClassification

  - Transition probabilities (e.g. *how likely is label I-PER, given the previous predicted label is B-PER*)

  - In theory, can produce invalid sequences (e.g. I-PER after O), but in practice does not really happen

- Adding a CRF (Conditional Random Field) on top of transformer adds direct knowledge of previous labels (dependencies between the predictions)

TURKUNLP
.ORG

UNIVERSITY
OF TURKU

# Sequence labeling with transformers

- Practical note: unlike most text classification datasets, sequence labeling datasets will normally come with a definition of what constitutes a "token" so a label (tags) can be assigned to each, e.g.
    - tokens = ["Se", "on", "Turun", "lähellä"]
    - labels = ["O", "O", "B-LOC", "O"]
- Transformer models have their own definitions of tokens, which cannot be readily changed. This adds a requirement to map from corpus tokens to transformer tokens and back.
- (The challenge here is mostly technical: the different definitions of "token" are mostly not an issue for machine learning model performance.)

TURKUNLP
.ORG

UNIVERSITY
OF TURKU

# Evaluation metrics

- Token-level classification accuracy (correct predictions out of all predictions) generally used to evaluate e.g. POS tagging

- For tasks involving marking spans (e.g. NER), performance typically measured on span level in terms of exact-match precision, recall and F1-score:

  - Compare predicted and gold standard spans in terms of (start-token, end-token, type)

  - Only triples where all values match between predicted and gold are correct

- Precision: fraction of predicted spans that are correct

- Recall: fraction of gold standard spans that are correctly predicted

- F1-score: balanced harmonic mean of precision and recall

# Model notebook

- https://github.com/TurkuNLP/textual-data-analysis-course/blob/main/sequence_labeling_example.ipynb

UNIVERSITY
OF TURKU

# Alternative approaches

- Fine-tuned decoder LMs for sequence labeling (token classification head)?

  - Causal attention mask not optimal for sequence labeling

  - *Nokia is a …*

- Like text classification, sequence labeling tasks can be approached also using generative models with in-context learning

  - Especially useful in cases where interested entities are not standard NER entities → No supervised training data

  - E.g. Hobbies, professions etc.

  - Generative models (language modelling head) read the whole text before generating an answer

UNIVERSITY OF TURKU

# Alternative approaches

- How to formulate a prompt for NER?

    - IOB-tagging prompt may not be the optimal

    - Ask the model to return / tag all entities of certain type from the text

        - "*List the person names appearing in this sentence …*"

        - "*Mark all person names appearing in this sentence with <PERSON> and </PERSON> tags.*"

    - Computationally heavier compared to multi-class approach, as this needs to be run separately for each entity type

        - Of course one can ask several types at the same time, but in general one-by-one seem to work better accuracy wise

# Summary

- Text classification maps a token sequence into a label (or set of labels), sequence labeling into one label per token

- Substantial number of pre-trained models and datasets for these tasks openly available

- Fine-tuned encoder LM still the standard method for sequence labeling

- In-context learning with generative models especially useful for entity types not annotated in standard NER datasets

TURKUNLP
.ORG

UNIVERSITY
OF TURKU