

# Part-of-speech tagging с использованием нейронных сетей

Даниил Анастасьев

Научный руководитель: Евгений Инденбом

Москва, 2018

- 1 Введение
- 2 Признаки
- 3 Функции потерь
- 4 Данные
- 5 Заключение

## Описание задачи

- Part-of-speech tagging — важный источник признаков для большинства NLP pipeline'ов.

## Описание задачи

- Part-of-speech tagging — важный источник признаков для большинства NLP pipeline'ов.
- Задача — найти грамматические значения (*теги*) всех слов в предложении:

У двери стоял                      **стол**                      секретарши , ...

NOUN

Animacy=Inan

Case=Nom

Gender=Masc

Number=Sing

The interviews took                      **place**                      two years ago .

NN

## Омонимичность тегов

Грамматическое значение слова почти невозможно определить, не принимая во внимание его контекст:

she	hated	lies
	VBD	
PRP	VBN	NNS
	JJ	VBZ

# Обзор поставленной задачи

- Чтобы решить любую задачу машинного обучения, нам нужны:

# Обзор поставленной задачи

- Чтобы решить любую задачу машинного обучения, нам нужны:
  - ❶ Модель;

# Обзор поставленной задачи

- Чтобы решить любую задачу машинного обучения, нам нужны:
  - 1 Модель;
  - 2 Данные для обучения;



# Обзор поставленной задачи

- Чтобы решить любую задачу машинного обучения, нам нужны:
  - 1 Модель;
  - 2 Данные для обучения;
  - 3 Признаки, извлекаемые из данных;

# Обзор поставленной задачи

- Чтобы решить любую задачу машинного обучения, нам нужны:
  - 1 Модель;
  - 2 Данные для обучения;
  - 3 Признаки, извлекаемые из данных;
  - 4 Функция потерь.

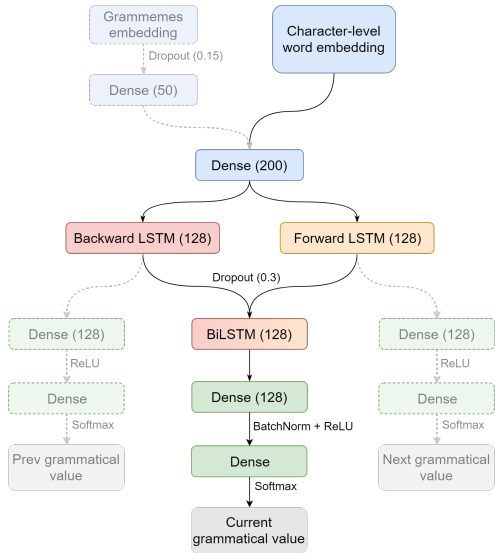
# Обзор поставленной задачи

- Чтобы решить любую задачу машинного обучения, нам нужны:
  - ❶ Модель;
  - ❷ Данные для обучения;
  - ❸ Признаки, извлекаемые из данных;
  - ❹ Функция потерь.
- Зафиксируем модель — Bidirectional LSTM;

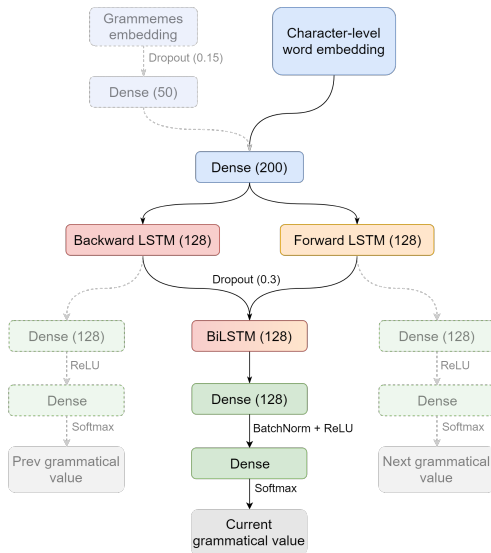
# Обзор поставленной задачи

- Чтобы решить любую задачу машинного обучения, нам нужны:
  - ❶ Модель;
  - ❷ Данные для обучения;
  - ❸ Признаки, извлекаемые из данных;
  - ❹ Функция потерь.
- Зафиксируем модель — Bidirectional LSTM;
- Улучшим её результат, сконцентрировавшись на остальных компонентах.

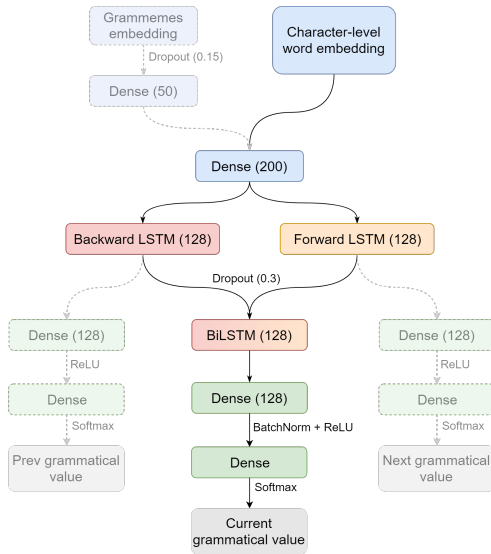
- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡



- Основой модели для нас послужит двуслойный BiLSTM;
- Эмбединги слов будет строить BiLSTM символического уровня — часть модели с самым высоким качеством на РТВ;



- Основой модели для нас послужит двуслойный BiLSTM;
- Эмбединги слов будет строить BiLSTM символического уровня — часть модели с самым высоким качеством на РТВ;
- Дополним пошагово данную модель новыми фишками.



# Датасеты

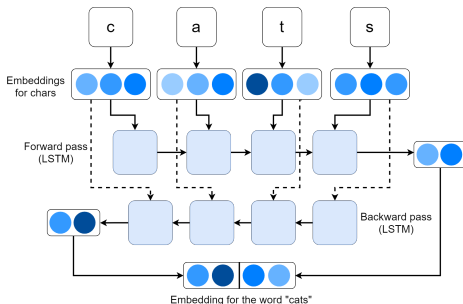
Использовались следующие датасеты для проверки моделей:

Датасет	Train	Dev	Test	#labels
PTB	912 344	131 768	129 654	45
UD SynTagRus	871 082	118 630	117 470	723
MorphoRuEval-2017	977 567	108 581	19 560	302
Tiger	711 041	88 152	89 054	54
UD Ukrainian IU	75 098	10 371	14 939	1 196



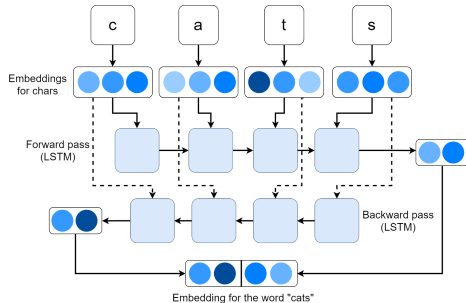
# Char BiLSTM

- BiLSTM символьного уровня — один из стандартных способов построить словный эмбединг;



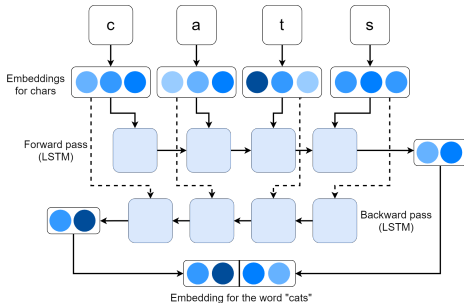
# Char BiLSTM

- BiLSTM символьного уровня — один из стандартных способов построить словный эмбединг;
- Обработывает символы один за другим;



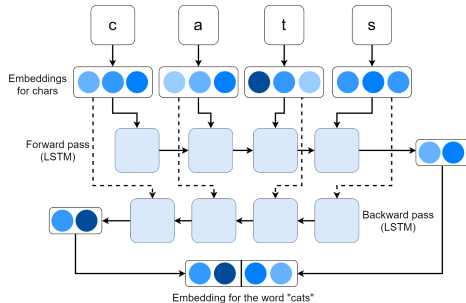
# Char BiLSTM

- BiLSTM символьного уровня — один из стандартных способов построить словный эмбединг;
- Обрабатывает символы один за другим;
- Может обрабатывать слова произвольной длины;



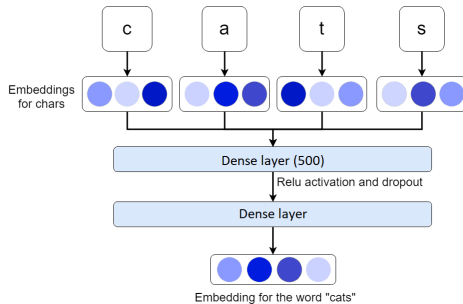
# Char BiLSTM

- BiLSTM символьного уровня — один из стандартных способов построить словный эмбединг;
- Обрабатывает символы один за другим;
- Может обрабатывать слова произвольной длины;
- Не параллелизуется эффективно.



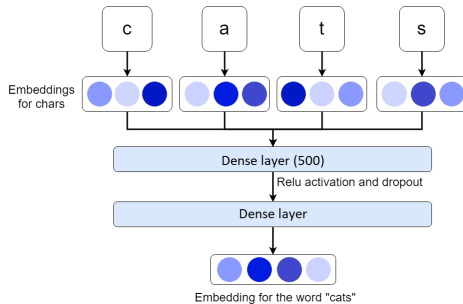
# Char FF

- Обычная feed-forward сеть — предлагаемая альтернатива Char BiLSTM;



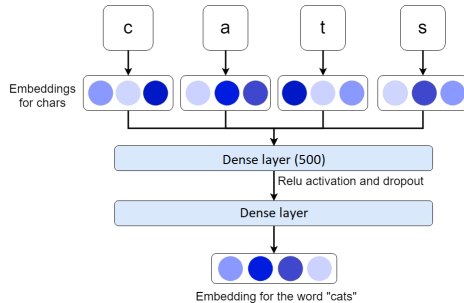
# Char FF

- Обычная feed-forward сеть — предлагаемая альтернатива Char BiLSTM;
- Обработывает конкатенацию символьных эмбедингов;



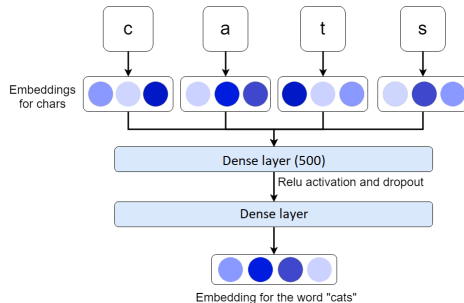
# Char FF

- Обычная feed-forward сеть — предлагаемая альтернатива Char BiLSTM;
- Обрабатывает конкатенацию символьных эмбедингов;
- Может работать лишь с фиксированной длиной слов:



# Char FF

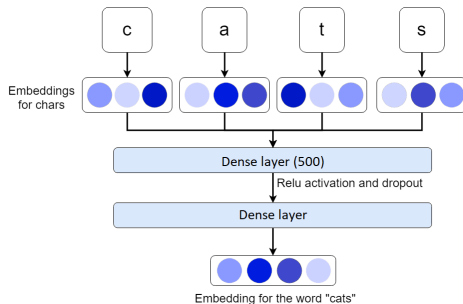
- Обычная feed-forward сеть — предлагаемая альтернатива Char BiLSTM;
- Обрабатывает конкатенацию символьных эмбедингов;
- Может работать лишь с фиксированной длиной слов:
  - 11–13 символов обычно достаточно;





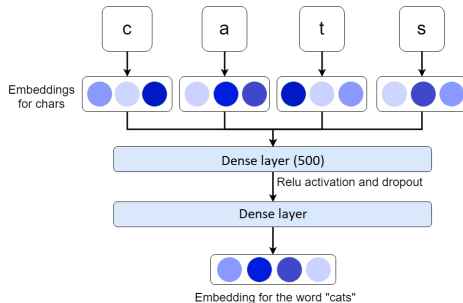
# Char FF

- Обычная feed-forward сеть — предлагаемая альтернатива Char BiLSTM;
- Обрабатывает конкатенацию символьных эмбедингов;
- Может работать лишь с фиксированной длиной слов:
  - 11–13 символов обычно достаточно;
  - Более короткие слова дополняются слева, более длинные обрезаются.



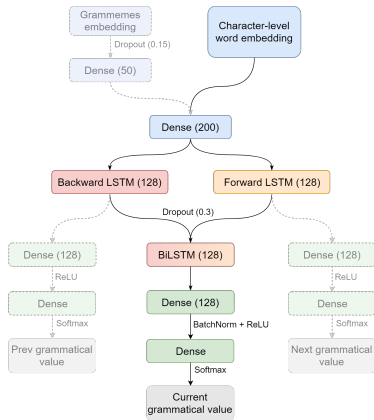
## Char FF

- Обычная feed-forward сеть — предлагаемая альтернатива Char BiLSTM;
- Обрабатывает конкатенацию символьных эмбедингов;
- Может работать лишь с фиксированной длиной слов:
  - 11–13 символов обычно достаточно;
  - Более короткие слова дополняются слева, более длинные обрезаются.
- Считаются гораздо быстрее Char BiLSTM.



# Сравнение вариантов эмбедингов

- По качеству модели близки;
- По скорости очевидно выигрывает Char FF.



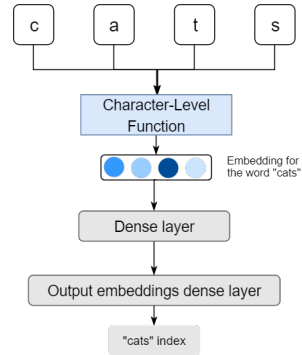
Dataset	Char BiLSTM	Char FF
PTB	97.02% / 96.98%	<b>97.32%</b> / <b>97.26%</b>
SynTagRus	<b>95.23%</b> / <b>95.39%</b>	94.98% / 95.16%
MorphoRuEval	96.48% / <b>94.69%</b>	<b>96.68%</b> / 94.63%
Tiger	98.27% / <b>99.86%</b>	<b>98.31%</b> / 99.73%
Ukrainian	80.10% / 78.70%	<b>81.51%</b> / <b>79.48%</b>

# Предобучение эмбедингов

- В предобученных словных эмбедингах закодирована важная информация, собранная на больших корпусах;

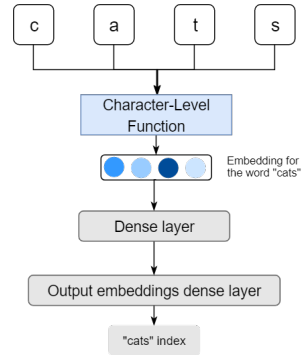
# Предобучение эмбедингов

- В предобученных словных эмбедингах закодирована важная информация, собранная на больших корпусах;
- Перенесём эту информацию с помощью автоэнкодера:



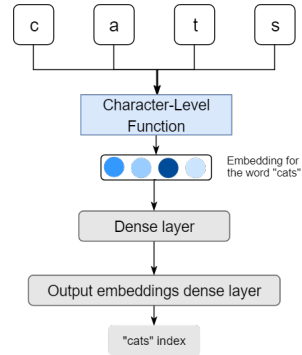
# Предобучение эмбедингов

- В предобученных словных эмбедингах закодирована важная информация, собранная на больших корпусах;
- Перенесём эту информацию с помощью автоэнкодера:
  - Энкодер — одна из функций символьного уровня;



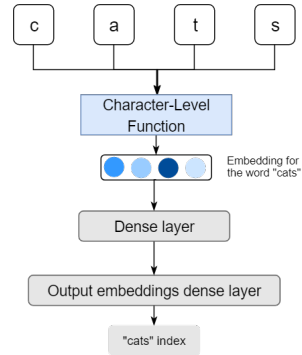
# Предобучение эмбедингов

- В предобученных словных эмбедингах закодирована важная информация, собранная на больших корпусах;
- Перенесём эту информацию с помощью автоэнкодера:
  - Энкодер — одна из функций символьного уровня;
  - Декодер — полносвязный слой, инициализированный предобученными эмбедингами.



# Предобучение эмбедингов

- В предобученных словных эмбедингах закодирована важная информация, собранная на больших корпусах;
- Перенесём эту информацию с помощью автоэнкодера:
  - Энкодер — одна из функций символьного уровня;
  - Декодер — полносвязный слой, инициализированный предобученными эмбедингами.
- Кросс-энтропийные потери стимулируют символьный эмбединг слова приближаться по косинусной мере к его предобученному варианту и удаляться от всех остальных эмбедингов.





## Результаты при применении предобучения

- Предобученные эмбединги обучались дальше вместе со всей моделью под задачу;

## Результаты при применении предобучения

- Предобученные эмбединги обучались дальше вместе со всей моделью под задачу;
- На первых эпохах модель с предобучением достигала заметно более высокого качества.

## Результаты при применении предобучения

- Предобученные эмбединги обучались дальше вместе со всей моделью под задачу;
- На первых эпохах модель с предобучением достигала заметно более высокого качества.

Dataset	Char FF	Char FF (Pretrained)
PTB	97.32% / 97.26%	<b>97.40%</b> / <b>97.31%</b>
SynTagRus	94.98% / 95.16%	<b>95.22%</b> / <b>95.36%</b>
MorphoRuEval	96.68% / 94.63%	<b>96.88%</b> / <b>94.63%</b>
Tiger	98.31% / <b>99.73%</b>	<b>98.39%</b> / 99.69%
Ukrainian	81.51% / 79.48%	<b>82.65%</b> / <b>80.67%</b>

# Граммемные эмбединги

- Сложно предсказать тег слова по набору его символов — близкие по написанию слова бывают очень далеки;

# Граммемные эмбединги

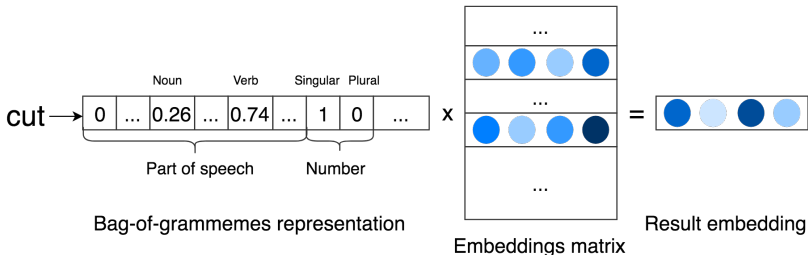
- Сложно предсказать тег слова по набору его символов — близкие по написанию слова бывают очень далеки;
- Будем оценивать априорные вероятности каждой из возможных граммем по словарю:

# Граммемные эмбединги

- Сложно предсказать тег слова по набору его символов — близкие по написанию слова бывают очень далеки;
- Будем оценивать априорные вероятности каждой из возможных граммем по словарю:
  - Например, форма существительного «cut» имеет частотность  $2.84 \cdot 10^{-5}$ , глагольная форма —  $8.75 \cdot 10^{-5}$ .  
Тогда  $P(\text{noun}) \approx 0.26$ .

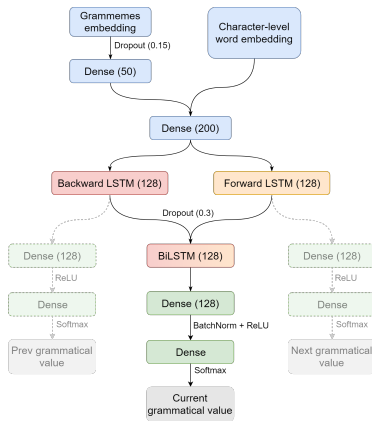
# Граммемные эмбединги

- Сложно предсказать тег слова по набору его символов — близкие по написанию слова бывают очень далеки;
- Будем оценивать априорные вероятности каждой из возможных граммем по словарю:
  - Например, форма существительного «cut» имеет частотность  $2.84 \cdot 10^{-5}$ , глагольная форма —  $8.75 \cdot 10^{-5}$ . Тогда  $P(\text{noun}) \approx 0.26$ .
- Добавим полносвязный слой для сокращения размерности эмбедингов:



# Результаты при применении граммемных эмбеддингов

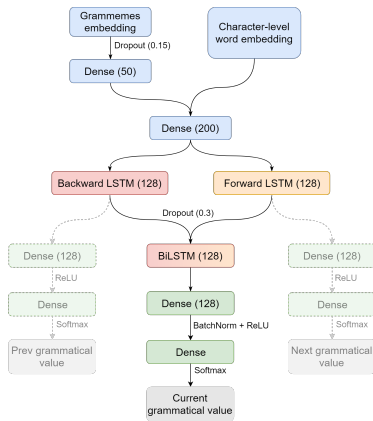
- На русском и украинском наиболее заметный прирост — до 35–43% ERR;
- На английском и немецком прирост незначительный.





# Результаты при применении граммемных эмбеддингов

- На русском и украинском наиболее заметный прирост — до 35–43% ERR;
- На английском и немецком прирост незначительный.



Dataset	Char FF (Pretrained)	+ Grammmemes
PTB	97.40% / <b>97.31%</b>	<b>97.43%</b> / 97.30%
SynTagRus	95.22% / 95.36%	<b>96.77%</b> / <b>97.00%</b>
MorphoRuEval	96.88% / 94.63%	<b>98.07%</b> / <b>95.36%</b>
Tiger	98.39% / 99.69%	<b>98.70%</b> / <b>99.85%</b>
Ukrainian	82.65% / 80.67%	<b>89.61%</b> / <b>88.06%</b>

# Результаты при добавлении языкового моделирования

- Добавим в модель потери от языкового моделирования;

# Результаты при добавлении языкового моделирования

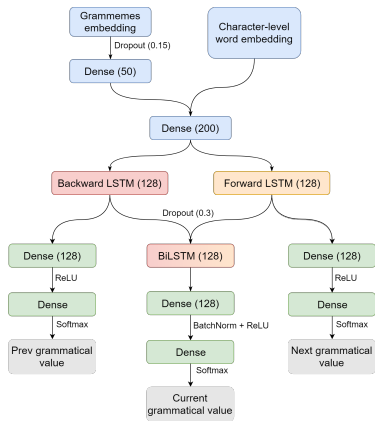
- Добавим в модель потери от языкового моделирования;
- *POS LM* пытается вместе с предсказанием тега слова выдавать теги предыдущего и следующего слов:

Forward LSTM(she, hated)  $\sim$   $\begin{matrix} \text{tag(hated)} \\ \text{tag(lies)} \end{matrix}$

# Результаты при добавлении языкового моделирования

- Добавим в модель потери от языкового моделирования;
- *POS LM* пытается вместе с предсказанием тега слова выдавать теги предыдущего и следующего слов:

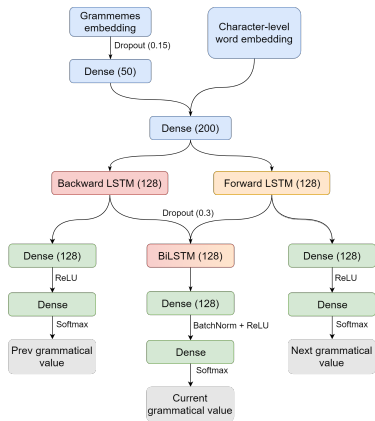
Forward LSTM(she, hated)  $\sim$   $\begin{matrix} \text{tag(hated)} \\ \text{tag(lies)} \end{matrix}$



# Результаты при добавлении языкового моделирования

- Добавим в модель потери от языкового моделирования;
- *POS LM* пытается вместе с предсказанием тега слова выдавать теги предыдущего и следующего слов:

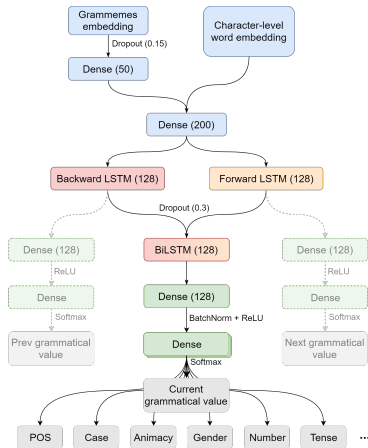
Forward LSTM(she, hated)  $\sim$  tag(hated)  
tag(lies)



Dataset	All features	+ POS LM
PTB	97.43% / 98.30%	<b>97.57%</b> / <b>97.49%</b>
SynTagRus	96.77% / 97.00%	<b>96.97%</b> / <b>97.24%</b>
MorphoRuEval	98.07% / 94.85%	<b>98.12%</b> / <b>96.72%</b>
Tiger	98.70% / <b>99.85%</b>	<b>98.71%</b> / 99.57%
Ukrainian	<b>89.61%</b> / 88.06%	89.48% / <b>88.07%</b>

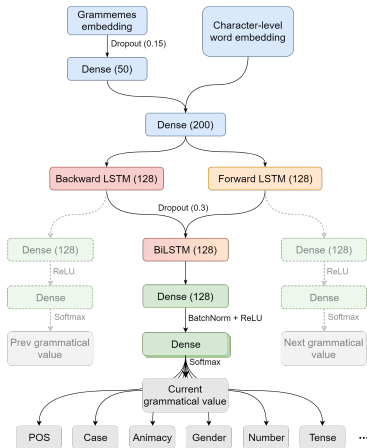
# Предсказание отдельных грамем

- Добавим в модель предсказание отдельных грамем для каждой из возможных грамматических категорий:



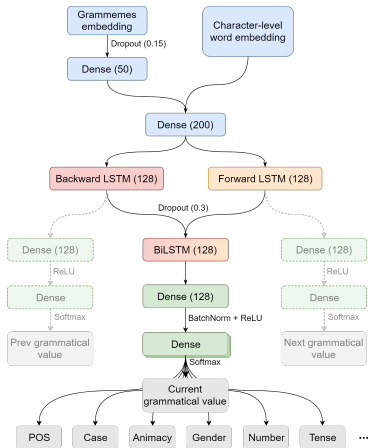
# Предсказание отдельных грамем

- Добавим в модель предсказание отдельных грамем для каждой из возможных грамматических категорий:
- Например, для категории NUMBER будем предсказывать распределение на граммах SINGLE, PLURAL и UNDEFINED.



# Предсказание отдельных грамем

- Добавим в модель предсказание отдельных грамем для каждой из возможных грамматических категорий:
- Например, для категории NUMBER будем предсказывать распределение на граммах SINGLE, PLURAL и UNDEFINED.



Dataset	All features	+ POS LM	+ Gram categories
SynTagRus	96.77% / 97.00%	<b>96.97% / 97.24%</b>	96.89% / 97.20%
MorphoRuEval	98.07% / 94.85%	<b>98.12% / 96.72%</b>	98.01% / 96.65%
Ukrainian	89.61% / 88.06%	89.48% / 88.07%	<b>90.17% / 89.01%</b>



# Перенос модели между датасетами

- Перенесем модель на SynTagRus с
  - 1 MorphoRuEval датасета с похожим UD тегсетом;
  - 2 Размеченным Comreno датасетом с не слишком похожим тегсетом.
- При переносе заменим выходной слой и будем первые несколько эпох тренировать только его, а только затем — всю модель целиком.

# Перенос модели между датасетами

- Перенесем модель на SynTagRus с
  - 1 MorphoRuEval датасета с похожим UD тегсетом;
  - 2 Размеченным Compreno датасетом с не слишком похожим тегсетом.
- При переносе заменим выходной слой и будем первые несколько эпох тренировать только его, а только затем — всю модель целиком.

Model	Accuracy
Best previous	96.97% / 97.24%
MorphoRuEval pretrained	<b>98.21%</b> / <b>98.33%</b>
Compreno pretrained	98.18% / 98.29%

# Перенос модели между языками

- Перенесем модель на украинский язык с SynTagRus аналогичным образом;

# Перенос модели между языками

- Перенесем модель на украинский язык с SynTagRus аналогичным образом;
- Для этого объединим все граммемы и символы, встречающиеся в этих языках.

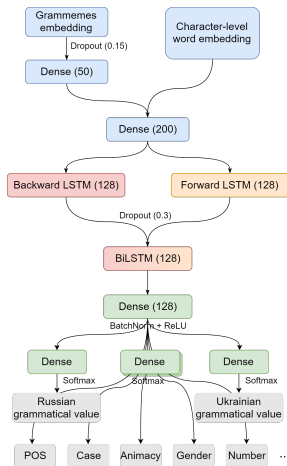
Dataset	Char FF + Grammemes + PosLM	+ Pretrained on Syntagrus
Ukrainian	89.48% / 88.07%	90.93% / 89.54%

# Совместная тренировка модели под несколько языков

- Модель при переносе очень быстро забывает то, на чём она училась до этого;

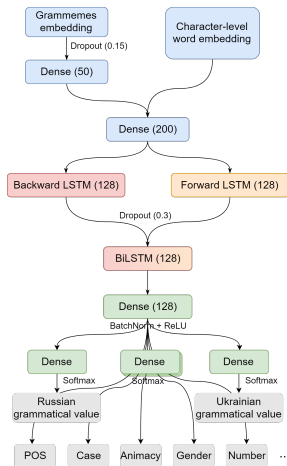
# Совместная тренировка модели под несколько языков

- Модель при переносе очень быстро забывает то, на чём она училась до этого;
- Будем тренировать модель под два языка сразу;



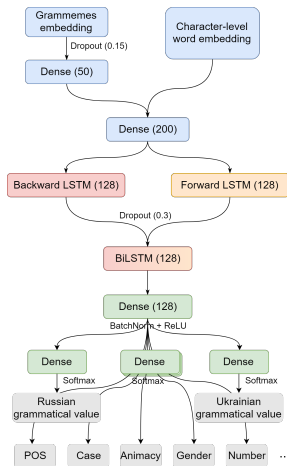
# Совместная тренировка модели под несколько языков

- Модель при переносе очень быстро забывает то, на чём она училась до этого;
- Будем тренировать модель под два языка сразу;
- Добавим к модели предсказание отдельных грамммем;



# Совместная тренировка модели под несколько языков

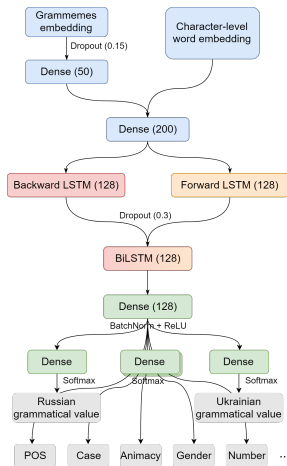
- Модель при переносе очень быстро забывает то, на чём она училась до этого;
- Будем тренировать модель под два языка сразу;
- Добавим к модели предсказание отдельных грамммем;
- Слои предсказания грамммем будут не language specific.





# Совместная тренировка модели под несколько языков

- Модель при переносе очень быстро забывает то, на чём она училась до этого;
- Будем тренировать модель под два языка сразу;
- Добавим к модели предсказание отдельных грамммем;
- Слои предсказания грамммем будут не language specific.



Dataset	Transfer baseline	Multi-lang	Multi-lang + Gram categories
Ukrainian	90.93% / 89.54%	91.15% / 89.33%	<b>91.72% / 89.87%</b>
Syntagrus	<b>96.77%</b> / 97.00%	96.44% / 96.69%	96.66% / <b>97.01%</b>

## Сравнение с baseline

- С использованием всех улучшений оказывается возможным значительно превзойти baseline;

Dataset	Char BiLSTM	Best Model	ERR
PTB	97.02% / 96.98%	<b>97.60%</b> / <b>97.51%</b>	19.4% / 17.5%
SynTagRus	95.23% / 95.39%	<b>98.21%</b> / <b>98.33%</b>	62.5% / 63.8%
MorphoRuEval	96.48% / 94.69%	<b>98.12%</b> / <b>96.72%</b>	46.5% / 38.2%
Tiger	98.27% / 99.86%	<b>98.74%</b> / <b>99.91%</b>	27.2% / 35.7%
Ukrainian	80.10% / 78.70%	<b>91.72%</b> / <b>89.87%</b>	58.4% / 52.4%

## Сравнение с baseline

- С использованием всех улучшений оказывается возможным значительно превзойти baseline;
- При этом размер модели почти не вырос;

Dataset	Char BiLSTM	Best Model	ERR
PTB	97.02% / 96.98%	<b>97.60%</b> / <b>97.51%</b>	19.4% / 17.5%
SynTagRus	95.23% / 95.39%	<b>98.21%</b> / <b>98.33%</b>	62.5% / 63.8%
MorphoRuEval	96.48% / 94.69%	<b>98.12%</b> / <b>96.72%</b>	46.5% / 38.2%
Tiger	98.27% / 99.86%	<b>98.74%</b> / <b>99.91%</b>	27.2% / 35.7%
Ukrainian	80.10% / 78.70%	<b>91.72%</b> / <b>89.87%</b>	58.4% / 52.4%

## Сравнение с baseline

- С использованием всех улучшений оказывается возможным значительно превзойти baseline;
- При этом размер модели почти не вырос;
- Итоговая модель значительно меньше большинства state-of-the-art моделей, но показывает сопоставимое качество.

Dataset	Char BiLSTM	Best Model	ERR
PTB	97.02% / 96.98%	<b>97.60%</b> / <b>97.51%</b>	19.4% / 17.5%
SynTagRus	95.23% / 95.39%	<b>98.21%</b> / <b>98.33%</b>	62.5% / 63.8%
MorphoRuEval	96.48% / 94.69%	<b>98.12%</b> / <b>96.72%</b>	46.5% / 38.2%
Tiger	98.27% / 99.86%	<b>98.74%</b> / <b>99.91%</b>	27.2% / 35.7%
Ukrainian	80.10% / 78.70%	<b>91.72%</b> / <b>89.87%</b>	58.4% / 52.4%

# Выводы

- Мы начали с сильного baseline — BiLSTM модель с Char BiLSTM эмбедингами;

# Выводы

- Мы начали с сильного baseline — BiLSTM модель с Char BiLSTM эмбедингами;
- Был предложен более быстрый аналог Char BiLSTM, показывающий сопоставимое качество;

# Выводы

- Мы начали с сильного baseline — BiLSTM модель с Char BiLSTM эмбедингами;
- Был предложен более быстрый аналог Char BiLSTM, показывающий сопоставимое качество;
- Был разработан метод для предобучения эмбедингов символьного уровня;

# Выводы

- Мы начали с сильного baseline — BiLSTM модель с Char BiLSTM эмбедингами;
- Был предложен более быстрый аналог Char BiLSTM, показывающий сопоставимое качество;
- Был разработан метод для предобучения эмбедингов символьного уровня;
- Были описаны дополнительные функции потерь, улучшающие качество целевой функции — POS LM и предсказание отдельных грамем;



# Выводы

- Мы начали с сильного baseline — BiLSTM модель с Char BiLSTM эмбедингами;
- Был предложен более быстрый аналог Char BiLSTM, показывающий сопоставимое качество;
- Был разработан метод для предобучения эмбедингов символьного уровня;
- Были описаны дополнительные функции потерь, улучшающие качество целевой функции — POS LM и предсказание отдельных грамем;
- Был продемонстрирован положительный эффект от переноса модели с датасета на датасет и с языка на язык;

## Выводы

- Мы начали с сильного baseline — BiLSTM модель с Char BiLSTM эмбедингами;
- Был предложен более быстрый аналог Char BiLSTM, показывающий сопоставимое качество;
- Был разработан метод для предобучения эмбедингов символьного уровня;
- Были описаны дополнительные функции потерь, улучшающие качество целевой функции — POS LM и предсказание отдельных грамем;
- Был продемонстрирован положительный эффект от переноса модели с датасета на датасет и с языка на язык;
- Был улучшен способ совместной тренировки модели на несколько языков;

## Выводы

- Мы начали с сильного baseline — BiLSTM модель с Char BiLSTM эмбедингами;
- Был предложен более быстрый аналог Char BiLSTM, показывающий сопоставимое качество;
- Был разработан метод для предобучения эмбедингов символьного уровня;
- Были описаны дополнительные функции потерь, улучшающие качество целевой функции — POS LM и предсказание отдельных грамем;
- Был продемонстрирован положительный эффект от переноса модели с датасета на датасет и с языка на язык;
- Был улучшен способ совместной тренировки модели на несколько языков;
- Всё это привело к значительному увеличению качества модели POS tagging'a.