

# **Сравнение стратегий подбора обучающих примеров для few-shot оценки языковой компетенции LLM** (на примере датасета BLiMP и модели Mistral)

---

Проектные задачи компьютерной лингвистики  
Элеонора Измайлова  
24.02.2025

# 1. Датасет

---

- **BLiMP**: The Benchmark of Linguistic Minimal Pairs for English [Warstadt et al., TACL 2020]
- создавался для выявления сильных и слабых сторон моделей (*n-gram*, *LSTM*, *Transformer-XL*, *GPT-2*) по сравнению с человеческими оценками
- состоит из **67** подразделов ('парадигм'), каждый из которых содержит **1 000** минимальных пар предложений на одно грамматическое явление
- был сгенерирован **автоматически** с сохранением структурной аналогии до ключевой позиции контраста и использованием словаря с многоуровневой разметкой размером **3000** слов

Phenomenon	N	Acceptable Example	Unacceptable Example
ANAPHOR AGR.	2	<i>Many girls insulted <u>themselves</u>.</i>	<i>Many girls insulted <u>herself</u>.</i>
ARG. STRUCTURE	9	<i>Rose wasn't <u>disturbing</u> Mark.</i>	<i>Rose wasn't <u>boasting</u> Mark.</i>
BINDING	7	<i>Carlos said that Lori helped <u>him</u>.</i>	<i>Carlos said that Lori helped <u>himself</u>.</i>
CONTROL/RAISING	5	<i>There was <u>bound</u> to be a fish escaping.</i>	<i>There was <u>unable</u> to be a fish escaping.</i>
DET.-NOUN AGR.	8	<i>Rachelle had bought that <u>chair</u>.</i>	<i>Rachelle had bought that <u>chairs</u>.</i>
ELLIPSIS	2	<i>Anne's doctor cleans one <u>important</u> book and Stacey cleans a few.</i>	<i>Anne's doctor cleans one book and Stacey cleans a few <u>important</u>.</i>
FILLER-GAP	7	<i>Brett knew <u>what</u> many waiters find.</i>	<i>Brett knew <u>that</u> many waiters find.</i>
IRREGULAR FORMS	2	<i>Aaron <u>broke</u> the unicycle.</i>	<i>Aaron <u>broken</u> the unicycle.</i>
ISLAND EFFECTS	8	<i>Whose <u>hat</u> should Tonya wear?</i>	<i>Whose should Tonya wear <u>hat</u>?</i>
NPI LICENSING	7	<i>The truck has <u>clearly</u> tipped over.</i>	<i>The truck has <u>ever</u> tipped over.</i>
QUANTIFIERS	4	<i>No boy knew <u>fewer than</u> six guys.</i>	<i>No boy knew <u>at most</u> six guys.</i>
SUBJECT-VERB AGR.	6	<i>These casseroles <u>disgust</u> Kayla.</i>	<i>These casseroles <u>disgusts</u> Kayla.</i>

Table 2: Minimal pairs from each of the twelve linguistic phenomenon categories covered by BLiMP. Differences are underlined. *N* is the number of 1,000-example minimal pair paradigms within each broad category.

## 2. 'Предыдущие' исследования

Тестирование в [Warstadt et al., TACL 2020]: «forced-choice» (модель сравнивает вероятности двух вариантов)

- минимальные пары создавались таким образом, чтобы различия в вероятностях объяснялись исключительно грамматическим контрастом, а не контекстом или длиной предложений

Model	Overall	ANA. AGR	ARG. STR	BINDING	CTRL. RAIS.	D-N AGR	ELLIPSIS	FILLER. GAP	IRREGULAR	ISLAND	NPI	QUANTIFIERS	S-V AGR
5-gram	60.5	47.9	71.9	64.4	68.5	70.0	36.9	58.1	79.5	53.7	45.5	53.5	60.3
LSTM	68.9	91.7	73.2	73.5	67.0	85.4	67.6	72.5	89.1	42.9	51.7	64.5	80.1
TXL	68.7	94.1	69.5	74.7	71.5	83.0	77.2	64.9	78.2	45.8	55.2	69.3	76.0
GPT-2	80.1	99.6	78.3	80.1	80.5	93.3	86.6	79.0	84.1	63.1	78.9	71.3	89.0
Human	88.6	97.5	90.0	87.3	83.9	92.2	85.0	86.9	97.0	84.9	88.1	86.6	90.9

Table 3: Percentage accuracy of four baseline models and raw human performance on BLiMP using a forced-choice task. A random guessing baseline would achieve an accuracy of 50%.

### 3. Данное (квази-)исследование

---

- тестирование путем **промπτинга**
  - ◆ **zero-shot** [Асс. 83.58]
  - ◆ **zero-shot CoT** [Асс. 87.01]
- размер тестовой выборки: **670** предложений (**10** предложений на каждое из **67** грамматических явлений)
- **Цель:** определить лучшую стратегию подбора обучающих примером для ***few-shot*** оценки

### 3. Данное (квази-)исследование

---

#### Эксперименты:

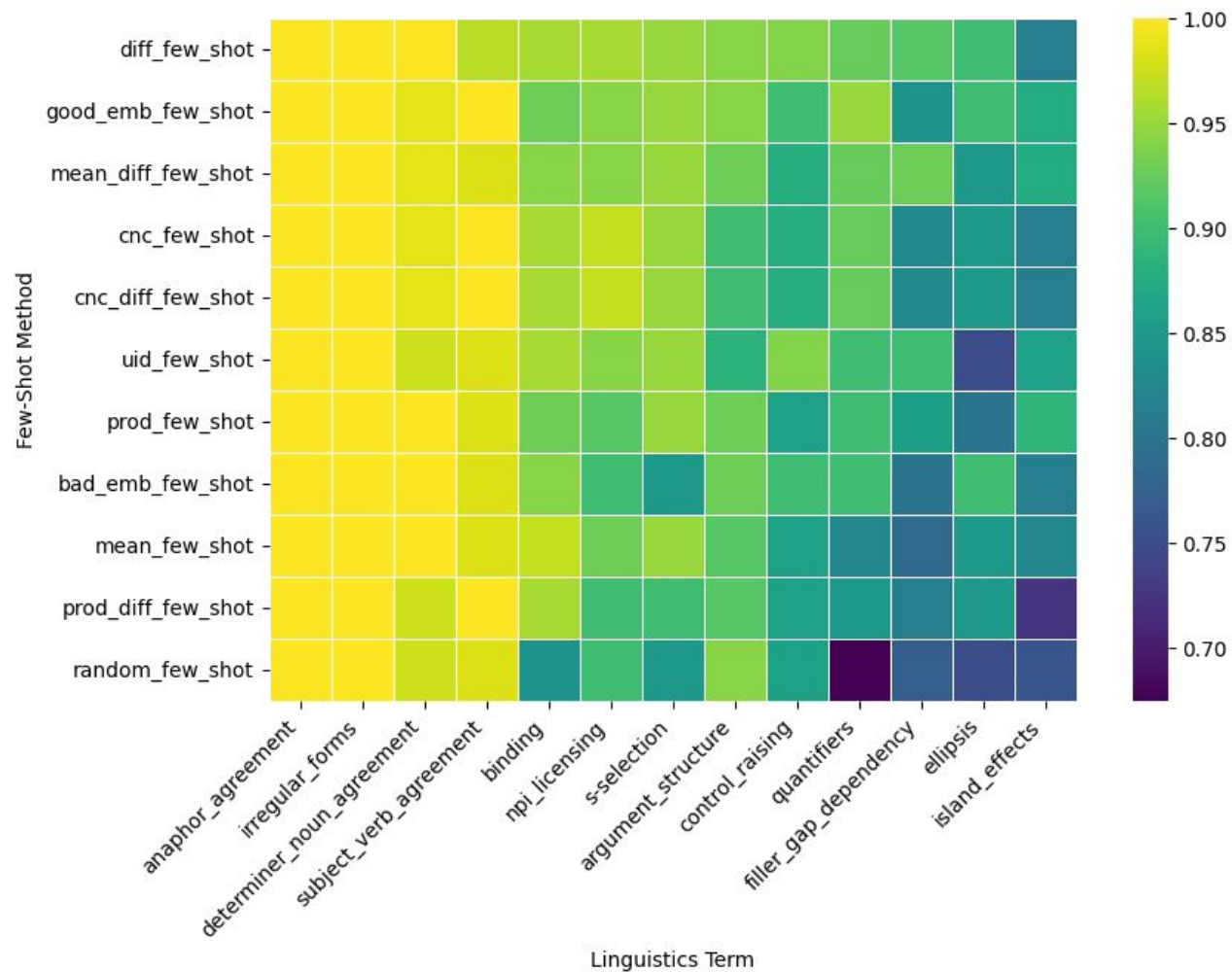
1. **Без использования NNS по эмбедингам**
  - a. случайный выбор из **всего** датасета
  - b. случайный выбор из **подраздела** датасета на то же грамматическое явление
2. **С использованием NNS (по эмбедингам предложений)**
  - a. поиск по близким к **грамматичному** предложению в тестовой паре
  - b. поиск по близким к **неграмматичному** предложению в тестовой паре
3. **С использованием NNS (по эмбедингам пар предложений)**
  - a. поиск по близким к **разности** эмбедингов тестовой пары
  - b. поиск по близким к **покомпонентному умножению** эмбедингов тестовой пары
  - c. поиск по близким к **усредненному** эмбедингу тестовой пары
  - d. поиск по близким к **конкатенированному** эмбедингу тестовой пары
4. **С использованием NNS (комбинированный подход)**
  - a. поиск по близким к конкатенации разности и **усреднения**
  - b. поиск по близким к конкатенации разности и **покомпонентного умножения**
  - c. поиск по близким к конкатенации разности и **конкатенации**

### 3. Данное (квази-)исследование

---

#### Эксперименты:

1. **Без использования NNS по эмбедингам**
  - a. случайный выбор из **всего** датасета [Асс. 86.87]
  - b. случайный выбор из **подраздела** датасета на то же грамматическое явление [Асс. 92.69]
2. **С использованием NNS (по эмбедингам предложений)**
  - a. поиск по близким к **грамматичному** предложению в тестовой паре [Асс. 93.28]
  - b. поиск по близким к **неграмматичному** предложению в тестовой паре [Асс. 91.04]
3. **С использованием NNS (по эмбедингам пар предложений)**
  - a. поиск по близким к **разности** эмбедингов тестовой пары [Асс. 93.73]
  - b. поиск по близким к **покомпонентному умножению** эмбедингов тестовой пары [Асс. 92.24]
  - c. поиск по близким к **усредненному** эмбедингу тестовой пары [Асс. 90.09]
  - d. поиск по близким к **конкатенированному** эмбедингу тестовой пары [Асс. 92.09]
4. **С использованием NNS (комбинированный подход)**
  - a. поиск по близким к конкатенации разности и **усреднения** [Асс. 93.58]
  - b. поиск по близким к конкатенации разности и **покомпонентного умножения** [Асс. 89.94]
  - c. поиск по близким к конкатенации разности и **конкатенации** [Асс. 91.94]





## 4. Дальнейшие планы

---

- сравнение с NNS по нормализованным эмбедингам
- дообучение матрицы проекций эмбедингов + NNS
- *TBD...*

## 5. **Использованная литература**

- ★ Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. [BLiMP: The Benchmark of Linguistic Minimal Pairs for English](#). *Transactions of the Association for Computational Linguistics*, 8:377–392.