

**ShiftRu**

**СеМаНТИчЕсКиЕ сдВИГи в  
РуСсКиХ нОвОсТяХ**

**Настя Чижикова и Женя Егорова**

# Проект ShiftRy

- ❖ ShiftRy 2020: анализ диахронических изменений в употреблении слов на российских новостях;
- ❖ период с 2010 по 2020;
- ❖ новостные веб-сайты газет с разными политическими взглядами: лояльные, оппозиционные, смешанные

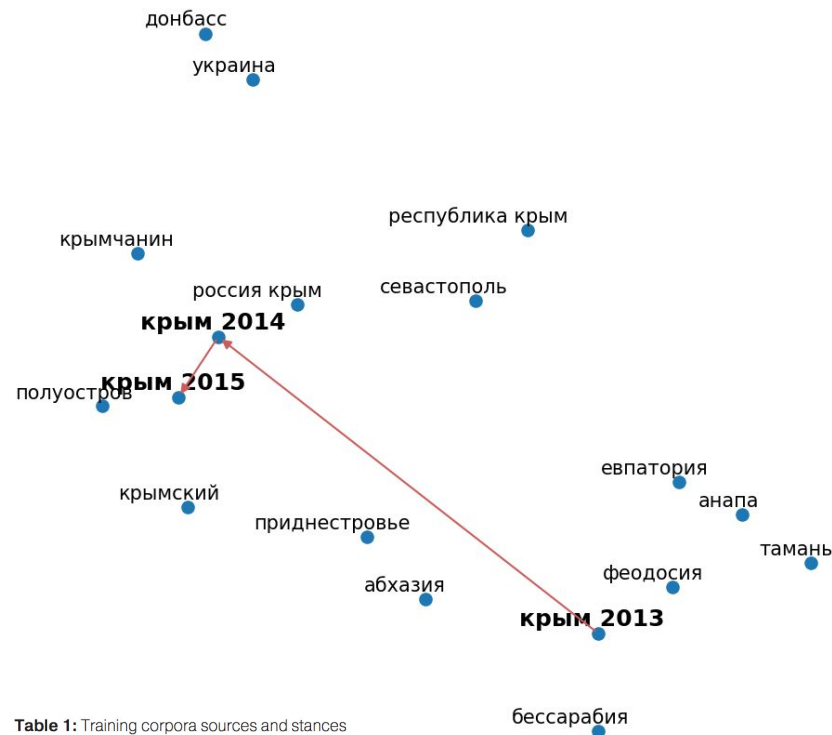


Table 1: Training corpora sources and stances

Nr.	Title	URL	Stance
1	Fontanka.ru	<a href="https://www.fontanka.ru/">https://www.fontanka.ru/</a>	Opposition
2	Gazeta.ru	<a href="https://www.gazeta.ru/">https://www.gazeta.ru/</a>	Loyal
3	Interfax	<a href="https://www.interfax.ru/">https://www.interfax.ru/</a>	Neutral
4	Izvestia	<a href="https://iz.ru/">https://iz.ru/</a>	Loyal
5	KP	<a href="https://www.kp.ru/">https://www.kp.ru/</a>	Loyal
6	Lenta.ru	<a href="https://lenta.ru/">https://lenta.ru/</a>	Mixed
7	Novaya Gazeta	<a href="https://novayagazeta.ru/">https://novayagazeta.ru/</a>	Opposition
8	N + 1	<a href="https://nplus1.ru/">https://nplus1.ru/</a>	Scientific
9	RBC	<a href="https://www.rbc.ru/">https://www.rbc.ru/</a>	Neutral
10	The Village	<a href="https://www.the-village.ru/">https://www.the-village.ru/</a>	Opposition

## Наша идея

- ❖ Рассмотреть отдельно провластные и оппозиционные СМИ, сохранив разделение по годам, и сравнить расположения векторов слов в пространстве;
- ❖ Попробовать получать ближайших соседей при помощи маскирующих языковых моделей и сравнить результаты по годам и политической направленности.

# Данные

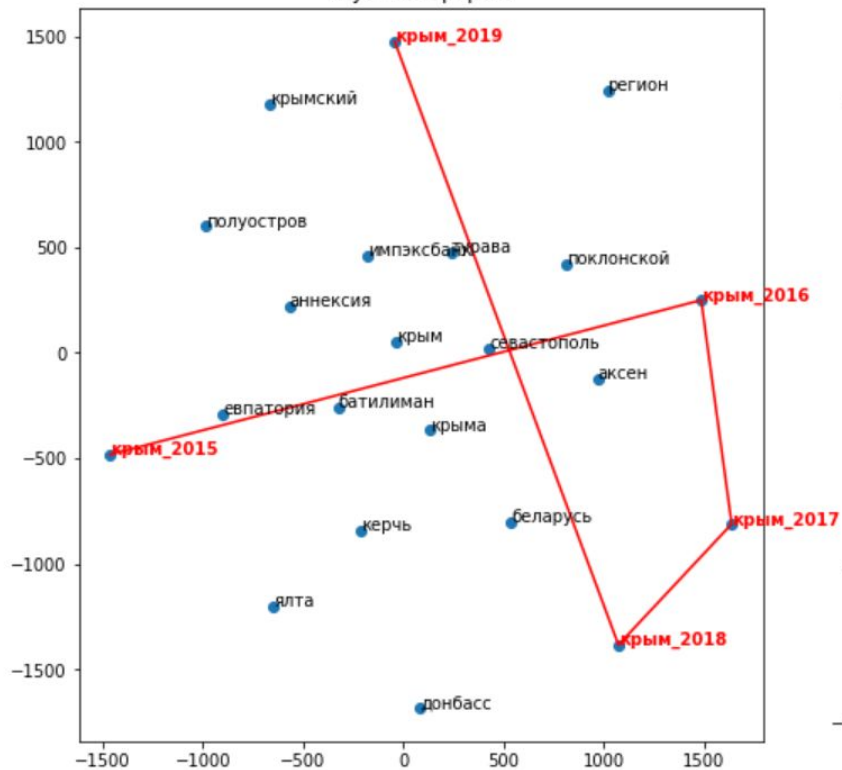
- данные с архива RusVectores с 2015 по 2019 год
- для каждого года – отдельный .txt файл для оппозиции и лоял  
(прошедший препроцессинг: токенизация (razdel), удаление пунктуации  
и стоп-слов, лемматизация и pos-tagging (nltk))

# Реализация

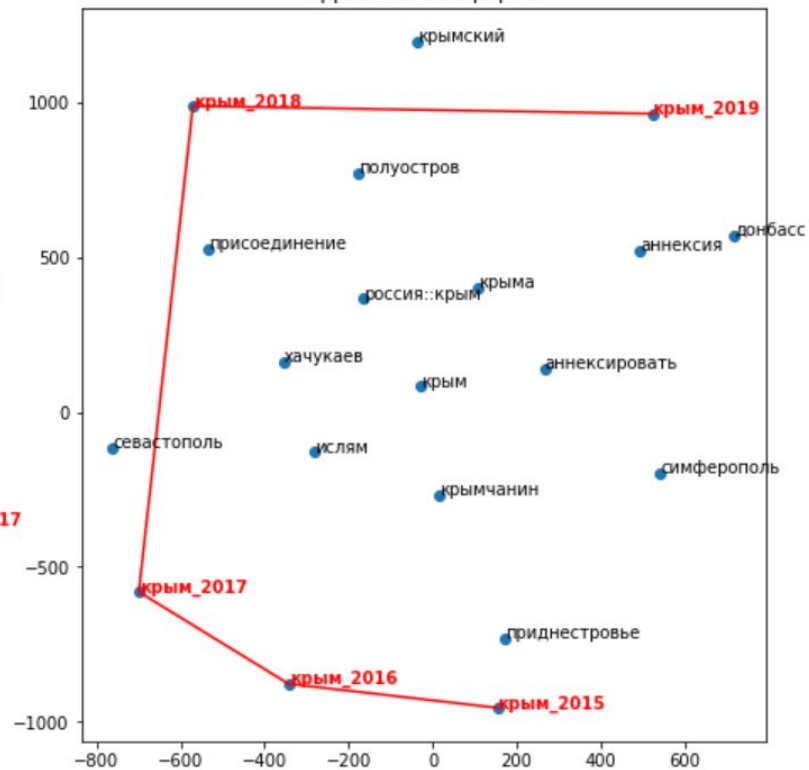
- ❖ за основу бралась модель с [RusVectors](#) 'taiga\_upos\_skipgram\_300\_2\_2018';
- ❖ данная модель дообучалась на новостных текстах каждого года, разделенных по политической направленности, полученные модели сохранялись;
- ❖ для построения векторов слов в одном пространстве модели выравнивались при помощи Procrustes transformation: алгоритм приближения матрицы  $X$  к матрице  $Y$  путем набора трансформаций
- ❖ при помощи t-SNE строились векторы заданного слова и 6 его ближайших соседей

# Результаты

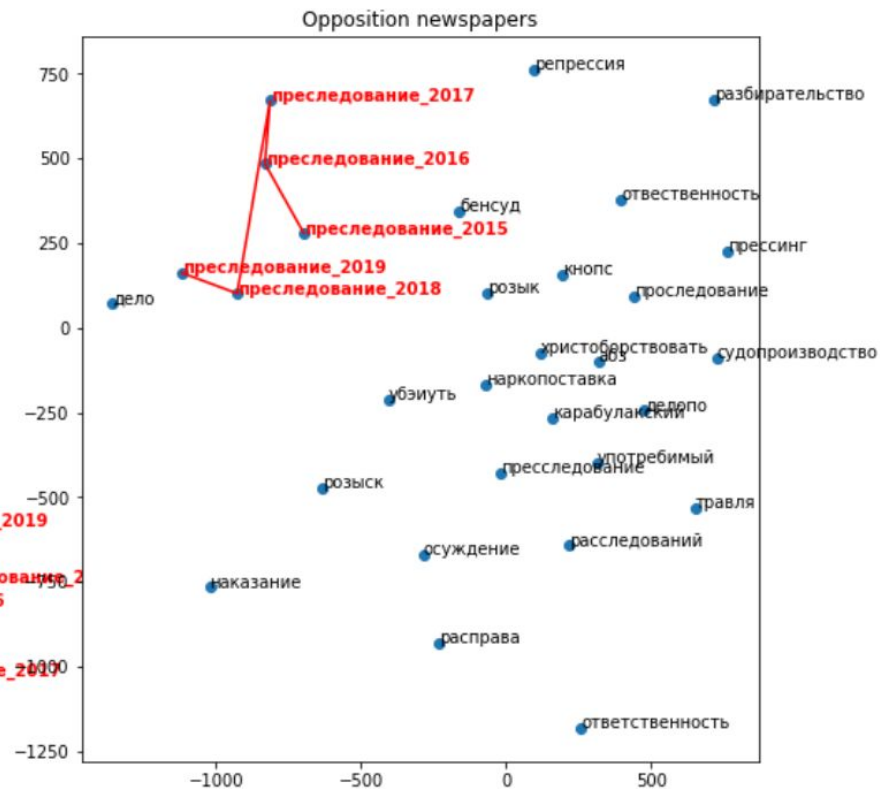
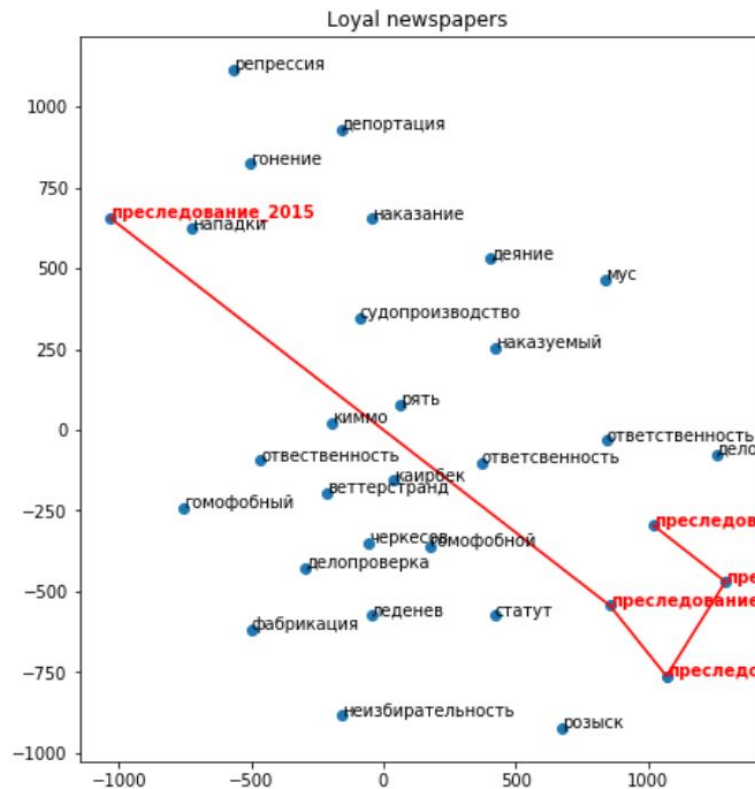
Loyal newspapers



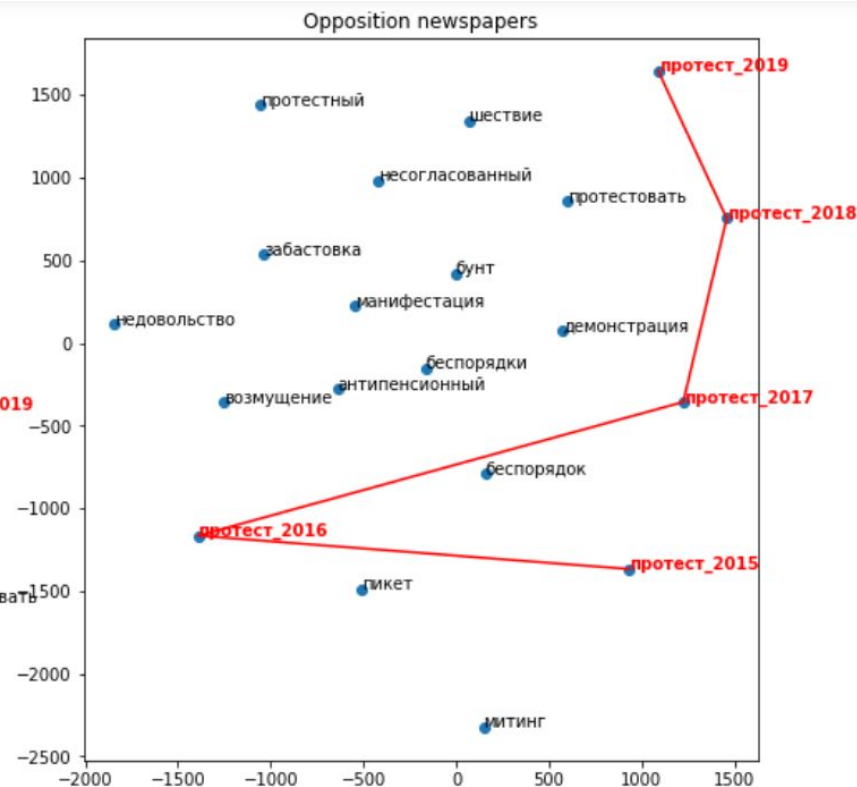
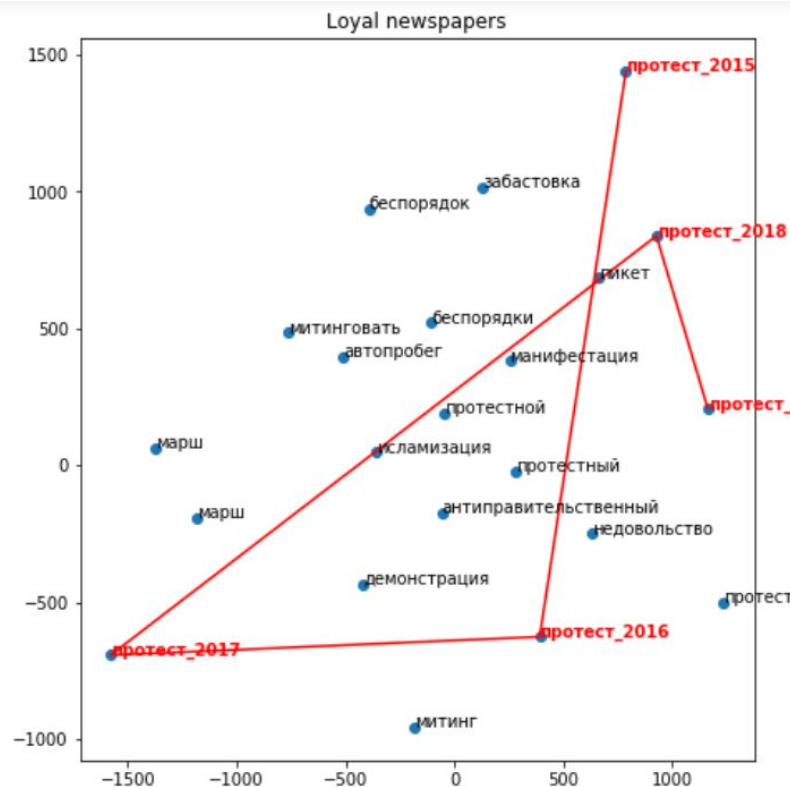
Opposition newspapers



# Результаты

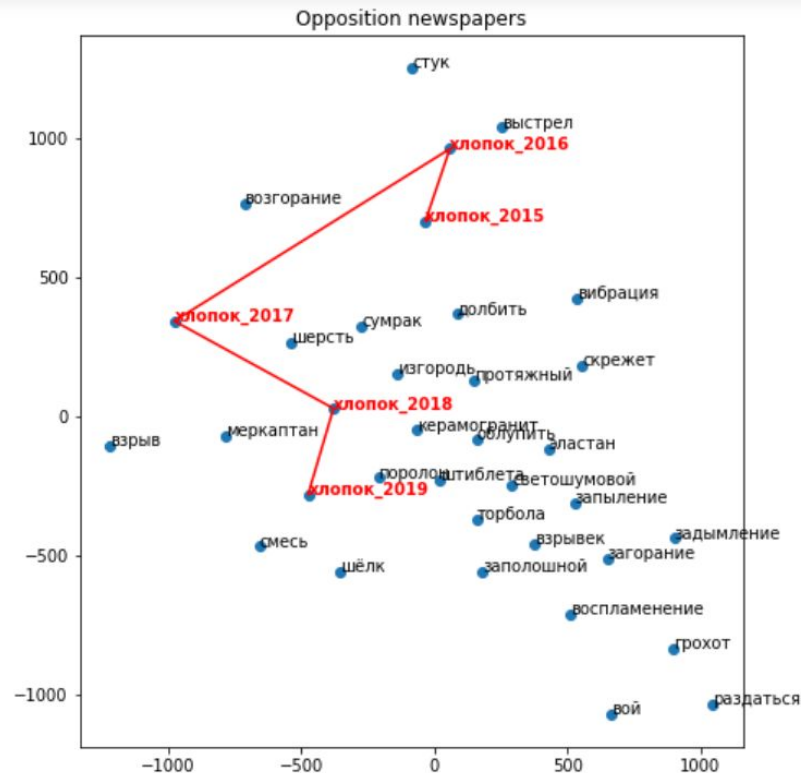
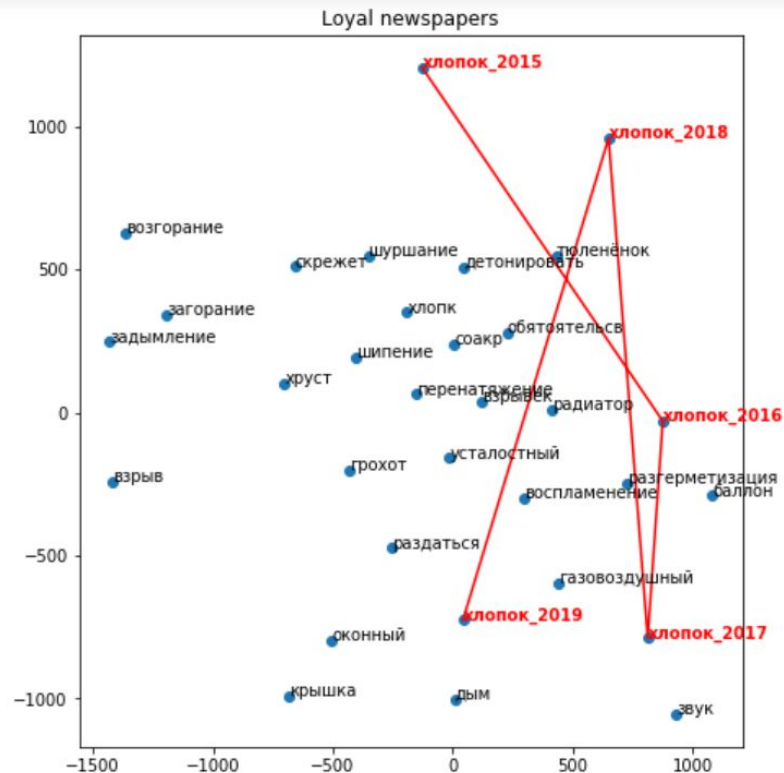


# Результаты





# Результаты



# Языковые модели?

- ❖ за основу брался русскоязычный BERT: `'rubert_cased_L-12_H-768_A-12_pt'`;
- ❖ в качестве данных были взяты те же тексты с убранными pos-тегами;
- ❖ модели дообучалась на восстановление маскированных токенов.

# Языковые модели

## "взрыв в [MASK] в Санкт-Петербурге"

2015	2016	2017	2018
['понедельник', 'год', 'четверг', 'день', 'воскресение', 'период']	['случиться', '##ити', '##еть', '[PAD]', '##еться', 'метро']	[' <b>метро</b> ', 'порту', 'понедельник', 'четверг', 'дом', 'спорткомплекс']	['четверг', 'час', 'понедельник', 'случиться', '##ити', 'пропасть']

## "Присоединение [MASK] к Российской Федерации"

2015	2016	2017	2018
[' <b>крым</b> ', 'полуостров', 'регион', 'россия', 'республика', 'территория']	['россия', 'страна', 'регион', 'это', 'федерация', 'государство']	['россия', 'крым', 'регион', 'страна', 'полуостров', 'киев']	['россия', 'страна', 'киев', 'крым', 'присоединение', 'республика']

взрыв в метро санкт-петербурге

Все

Картинки

Новости

Видео

Карты

Ещё

Настройки

Инструменты

Результатов: примерно 928 000 (0,64 сек.)

[https://ru.wikipedia.org/wiki/Теракт\\_в\\_Петербурге](https://ru.wikipedia.org/wiki/Теракт_в_Петербурге)

Теракт в Петербургском метрополитене (2017) — Википедия

Террористический акт в Петербургском метрополитене — взрыв, произошедший в понедельник, 3 апреля 2017 года в 14:33 в Санкт-Петербурге на ...

Цель нападения: Пассажиры метро

Способ нападения: взрыв, теракт-само...

Дата: 3 апреля 2017 года; 14:33 (местное в...

Подозреваемые: Абдор Азимов

Список террористических... · (1996) — Википедия · Февраль 2004

# Языковые модели

"[MASK] выиграл выборы в америке"

2015	2016	2017	2018
['год', 'также', 'политик', 'это', 'однако', 'спортсмен']	['миллиардер', 'год', 'трамп', 'это', 'накануне', 'ранее']	['также', 'однако', 'затем', 'кроме', 'это', 'трамп']	['год', 'политик', 'октябрь', 'мужчина', 'трамп', 'ноябрь']

"В России продолжаются преследования [MASK]."

Провластные			Оппозиционные		
2015	2016	2017	2015	2016	2017
['человек', 'ребенок', 'преступник', 'медведь', 'автомобиль', 'россиянин']	['человек', 'преступник', 'мужчина', 'ребенок', 'подозревать', 'женщина']	['подозревать', 'коррупция', 'террорист', 'боевик', 'человек', 'россия']	['р', 'россия', 'ст', 'обвинять', 'подозревать', 'задержать']	['страна', 'человек', 'россия', 'иностранец', 'воина', 'запад']	['человек', 'геев', 'коррупция', 'гражданин', 'год', 'задержать']

# Возможные улучшения

- дополнить коллекцию текстов по годам до 2014 и после 2020 и обучить на них модели;
- в языковых моделях использовать нелемматизированные тексты:

"Против [MASK] завели уголовное дело."

The Village 2019 (оппозиционно)	Lenta.ru 2019 (провластно)
['него', 'политика', 'студента', 'них', 'художника', 'ученого']	['него', 'нее', 'мужчины', 'девушки', 'злоумышленника', 'задержанного']

"В России продолжаются преследования [MASK]."

The Village 2019 (оппозиционно)	Lenta.ru 2019 (провластно)
['оппозиционеров', 'россиян', 'политзаключенных', 'активистов', 'протестующих', 'заключенных']	['животных', 'биатлонистов', 'россиян', 'пилотов', 'геев', 'спортсменов']

СПАСИБО 🐱 ЗА  
ВНИМАНИЕ!