



Проект создания газетного Интернет-корпуса

Михаил Дубов, Никита Левицкий, Егор Моренко,
Екатерина Черняк, Артем Шаль, Андрей Шестаков

Научно-учебная группа «Методы анализа и визуализации текстов»
Национальный исследовательский университет «Высшая школа экономики»



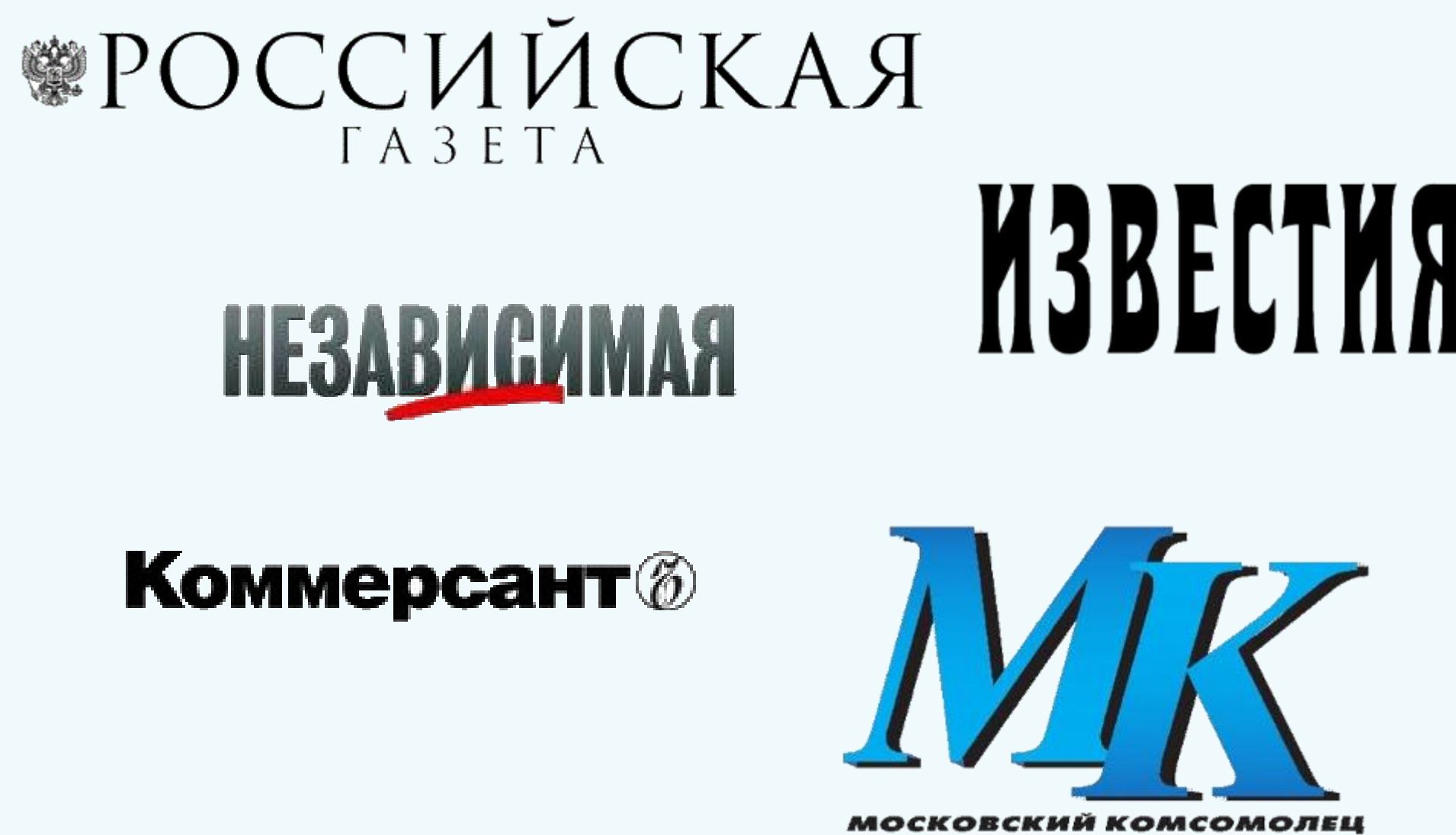
НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
УНИВЕРСИТЕТ

<http://ami.hse.ru/vitext>

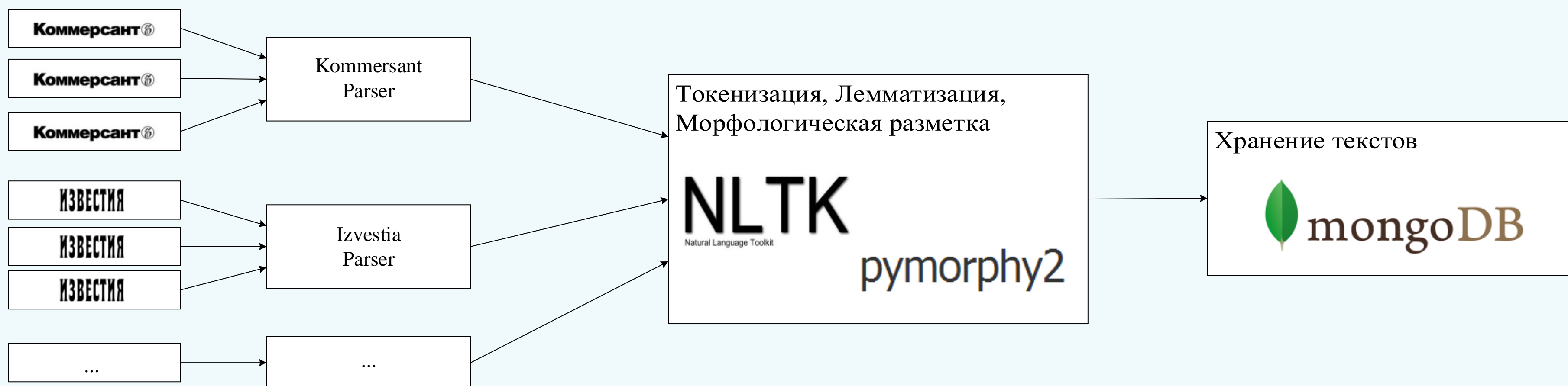
Почему газетные публикации?

- ▶ Тексты собираются не со всего веба;
- ▶ Фиксированное число источников:
 - ⇒ Качественное решение задачи удаления html-разметки и рекламы;
 - ⇒ Известная структура публикаций в каждом источнике.
- ▶ Высокое качество собранных данных:
 - Язык – достаточно формален, опечаток – мало;
 - Содержательно тексты – на высоком уровне;
 - Достаточно узкий набор тем;
 - Многие тексты категоризированы.

Источники текстов



Реализация

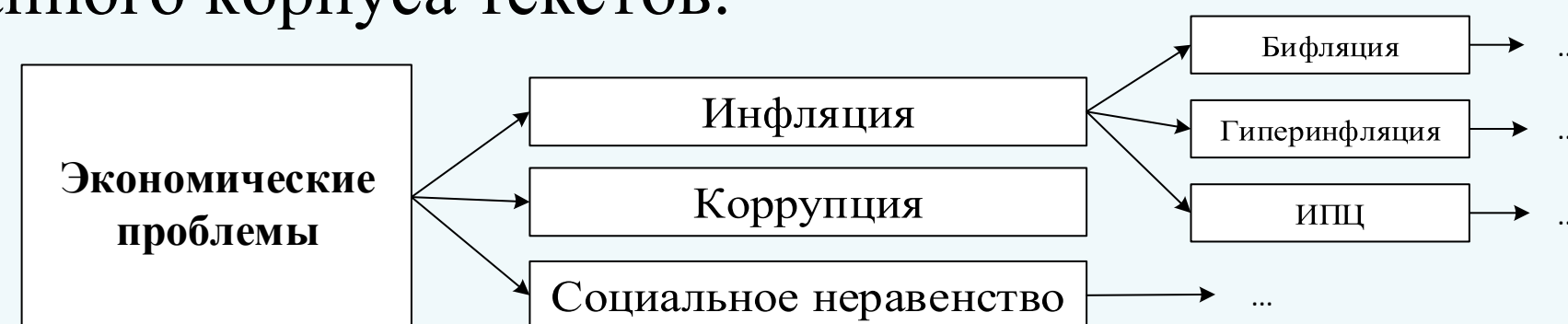


Основные принципы развития ПО для сбора корпуса

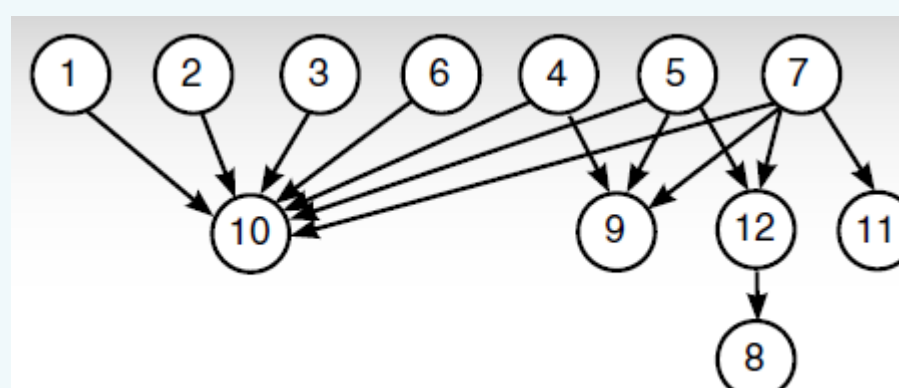
- ▶ Полная автоматизация процесса пополнения корпуса новыми текстами и их разметки;
- ▶ Простота расширяемости ПО парсерами для новых источников текстов;
- ▶ Возможность добавления в систему новых модулей по обработке текстов.

Анализ собранных текстов

- ▶ Построение таксономий и онтологий на основе данного корпуса текстов:



- ▶ Анализ связей между ключевыми словосочетаниями на данном корпусе текстов^[1]:



1. Ввод автоматизированного производства.
2. Выпуск пресс-релизов.
- ...
4. Изменение уровня концентрации собственности.
5. Повышение квалификации персонала.
- ...
9. Изменение организационно-правовой формы.
10. Повышение эффективности управления затратами.
- ...

^[1] Миркин Б. Г., Черняк Е. Л., Чугунова О. Н. Метод аннотированного суффиксного дерева для оценки степени вхождения строк в текстовые документы. Бизнес-информатика, №3(21), 2012. с. 31-41.