# Генератор парсеров

Кормышов Михаил

Руководитель: Фёдор Амосов

CSC, осень 2017

# Offline оценка качества поиска



Выбрать запросы → Скачать SEPR'а → Распарсить → Оценить → Посчитать метрики
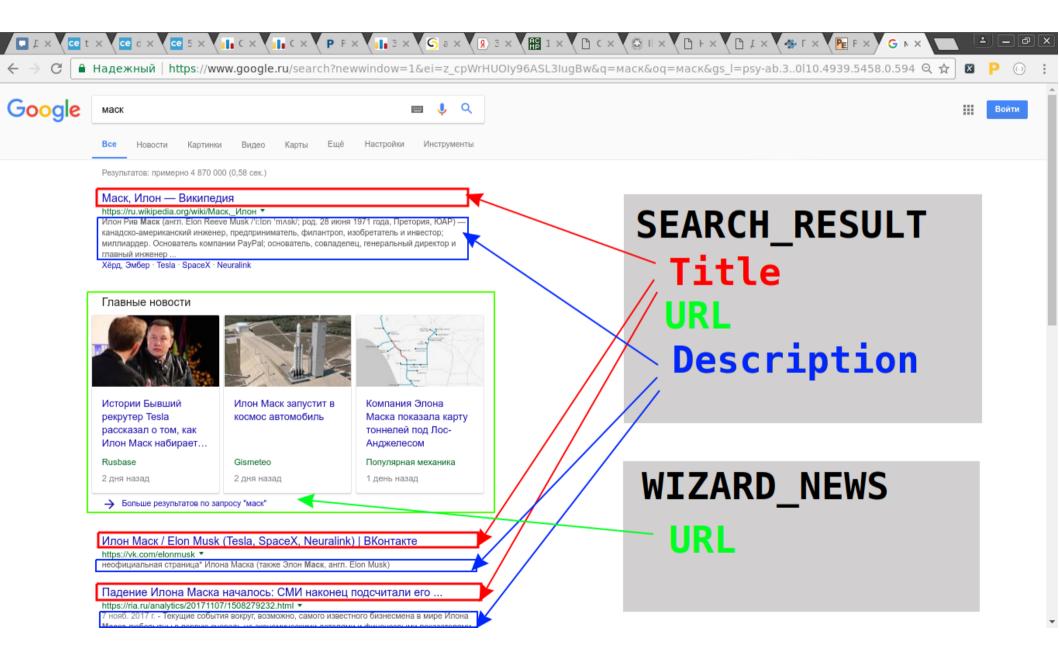
# Что парсим

# Как парсим сейчас

```python
            wizard.type = "WIZARD"
            wizard.wizard_type = "WIZARD_NEWS"
            wizard.alignment = "LEFT"
            wizard.page_url = self.get_from_page(element, ".", "href")
            wizard.title = self.get_from_page(element, ".", "string")
            return wizard

    def extract_markup(self, file_name):
        with open(file_name, "r") as file:
            tree = html.document_fromstring(file.read())
        markup = SearchMarkup()
        markup.file = file_name.split('/')[-1]
        markup.type = "HTMLTree"
        block_list = tree.xpath("//html/body/table/tbody/tr/td/div/div")
        for block in block_list:
            adv_list = block.xpath(".")
            for adv in adv_list:
                subblock_list = adv.xpath("./ol/li")
                for subblock in subblock_list:
                    if len(subblock.xpath("./h3")) > 0:
                        result = self.extract_adv(subblock)
                        markup.add(result)
            subblock_list = block.xpath("./div/div/ol/div")
            for subblock in subblock_list:
                document_list = subblock.xpath(".")
                for document in document_list:
                    if len(document.xpath("./div/div/cite")) > 0:
                        result = self.extract_search_result(document)
                        markup.add(result)
                wizard_image_list = subblock.xpath(".")
                for wizard_image in wizard_image_list:
                    if len(wizard_image.xpath("./div[1]/a")) > 0:
                        result = self.extract_wizard_image(wizard_image)
                        markup.add(result)
                wizard_news_list = subblock.xpath("./div/h3/a")
                for wizard_news in wizard_news_list:
```
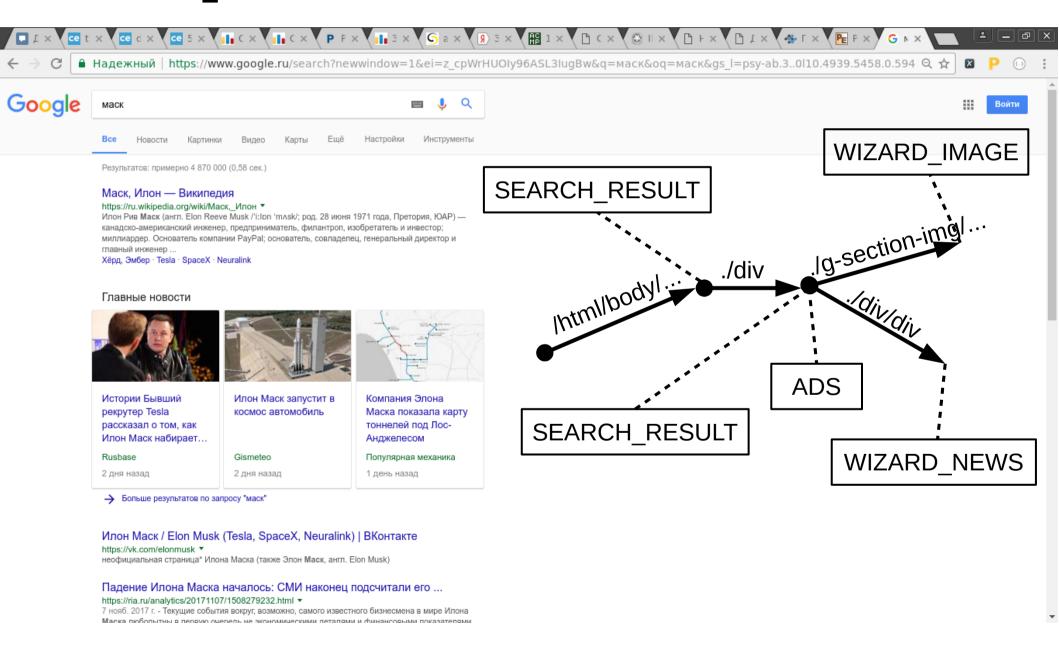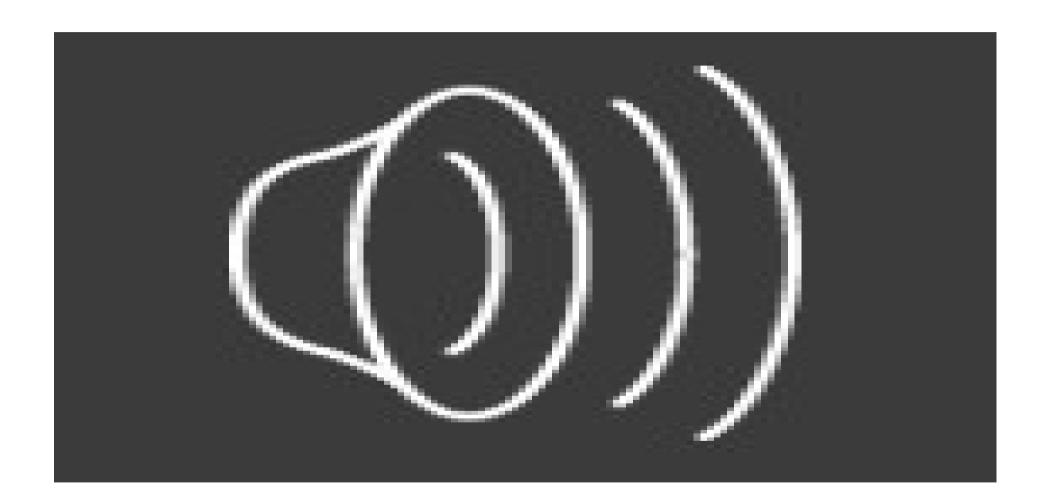
# Как хотим

# Аналоги

ScraperLab

A-Parser

# Алгоритм

# Технологии

- Python
- JavaScript
- XPath
- JSON
- CherryPy
- Chrome Extensions API

# Спасибо за внимание

Вопросы?