

## KWJP $\frac{1}{2}$ M

Paczka `samples.tar.gz` zawiera podkorpus KWJP $\frac{1}{2}$ M, czyli losowe próbki tekstów z Korpusu Współczesnego Języka Polskiego. KWJP $\frac{1}{2}$ M składa się ze 11 108 próbek obejmujących łącznie 513 009 słów tekstowych „od spacji do spacji”. Pojedyncza próbka to jedno lub więcej zdań występujących kolejno w tekście korpusu.

Liczba próbek pochodzących z każdego tytułu (książki, dziennika, tygodnika etc.) wchodzącego w skład KWJP jest proporcjonalna do udziału tego tytułu w stumilionowym korpusie zrównoważonym.

Próbki zapisane są w formacie `.json`. W pojedynczym pliku `xyz_sample.json` zebrane są wszystkie próbki wylosowane z dokumentu korpusu o identyfikatorze `xyz`. Przykładowa zawartość takiego pliku wygląda następująco:

```
{
  "meta": {
    "id": "146648",
    "meta_author": "",
    "meta_title": "Kącik Gier Historycznych: Ostrołęka 26 maja 1831",
    "meta_editor": "",
    "meta_main_title": "Rebel Times",
    "meta_published": "2015-11-24",
    "meta_first_published": "",
    "meta_publisher": "",
    "meta_pubplace": "Gdańsk",
    "meta_isbn": "",
    "meta_issn": "",
    "meta_pubnumber": "98",
    "meta_topic": "",
    "meta_kwjp_type": "typ_fakt",
    "meta_kwjp_channel": "kanal_prasa_miesiecznik",
    "meta_label": "rebel-times_98"
  },
  "samples": [
    {
      "sentence_ids": [
        "text_structure.xml/para_5/sent_4",
        "text_structure.xml/para_5/sent_5"
      ],
      "text": "Oddziały kawalerii są w skali dywizjonu (jeden dywizjon to dwa szwadrony, dwa dywizjony składają się zazwyczaj na pułk kawalerii), czyli jeden oddział to zwykle 200-400 ludzi. W przypadku artylerii oddziały są w skali baterii - jedna bateria to na ogół od 6 do 12 dział i 150-300 ludzi obsługi."
    },
    {
      "sentence_ids": [
        "text_structure.xml/para_18/sent_3",
        "text_structure.xml/para_18/sent_4",

```

```

        "text_structure.xml/para_18/sent_5"
    ],
    "text": "Trudno zrobić szybko dobrą, trochę bardziej zaawansowaną grę historyczną. Trzeba do tego i wiedzy, i umiejętności, i praktyki (trzeba przez dłuższy czas grać w takie gry). Trzeba mieć jakąś własną wizję, pomysł, jak taka gra ma wyglądać i czym będzie się różnić od innych jej podobnych."
  }
]
}

```

Słownik reprezentujący próbki z pojedynczego dokumentu zawiera dwa klucze: **meta** (metadane próbki) oraz **samples** (właściwe próbki).

Elementy (być może puste, jeśli dana informacja była nieznana lub niedostępna) metadanych to:

- **id** – identyfikator dokumentu,
- **meta\_author** – autor,
- **meta\_title** – tytuł tekstu,
- **meta\_editor** – redaktor,
- **meta\_main\_title** – tytuł główny (np. tytuł czasopisma),
- **meta\_published** – data wydania,
- **meta\_first\_published** – data pierwszego wydania,
- **meta\_publisher** – wydawca,
- **meta\_pubplace** – miejsce wydania,
- **meta\_isbn** – ISBN,
- **meta\_issn** – ISSN,
- **meta\_pubnumber** – numer periodyku (dotyczy tylko prasy),
- **meta\_topic** – tematyka (wartość opisowa; dotyczy tylko książek),
- **meta\_kwjp\_type** – typ; możliwe wartości:
  - **typ\_fikcja**,
  - **typ\_fakt**,
  - **typ\_publicystyka**,
- **meta\_kwjp\_channel** – kanał; możliwe wartości:
  - **kanal\_ksiazka**,
  - **kanal\_prasa\_dziennik**,
  - **kanal\_prasa\_tygodnik**,
  - **kanal\_prasa\_miesiecznik**,
  - **kanal\_prasa\_inne**,
  - **kanal\_internet**,
- **meta\_label** – etykieta.

W przypadku książek pole **meta\_title** zawiera tytuł książki, zaś pole **meta\_main\_title** pozostaje puste. W przypadku wydawnictw cyklicznych pole **meta\_title** zawiera tytuł dokumentu (np. artykułu), zaś pole **meta\_main\_title** – tytuł dziennika, tygodnika etc. W szczególności może być kilka plików o tym samym tytule głównym, zawierających próbki pochodzące z różnych dokumentów w obrębie jednego tytułu np. prasowego.

Próbki tekstów zapisane są w postaci listy słowników, które z kolei zawierają dwa klucze:

- **sentence\_ids** – lista identyfikatorów zdań składających się na próbkę w pełnym korpusie,
- **text** – tekst próbki.