# Information Retrievel with PostgreSQL

## Alexander Hebel

Heidelberg University
Institute of Computer Science
Database Systems Research Group
vx228@uni-heidelberg.de

Mai 6, 2020

# Outline

## Outline

# Task definition

- How looks and performs an IRS made of a relational database
- Similar to Apache Solr
- Finding different database models
- Python api for the database creation and communication
- Crawl Wikipages to gather some text data
- Special type in PostgreSQL named tsvector (full text search)

### First goal

Support some boolean search querys like AND

# Outline

1. Introduction

2. **Approach and realizations**

3. Custom C-functions in PostgreSQL

4. Rating sections vs. rating pages

5. Conclusion

## Realization

### Wiki crawler

- Based on package wikipedia version 1.4.0
- Takes number of pages and category as input
- Also searches in subcategories
- Variable level of subcategories

### Database pipeline

- Used package psycopg2 version 2.8.5
- custom converter for tsvector

**Introduction**
oo

**Realization**
ooo●o

**Custom C-functions**
ooo

**Rating comparison**
oooooo

**Conclusion**
ooo

# Database models

# Tsvector

## Possibilities

- Full text search
- GIN-Index
- Automatic tokenization and lemmatization
- Adding weights
- Predefined rating function

## Limitations

- The number of lexemes must be less than $2\hat{6}4$
- Max position value: 16383
- No more than 256 positions per lexeme
- Relative small set of manipulation methods
- Limited rating

### Example

{'a':1,6,10 'and':8 'cat':3 'fat':2,11 'mat':7 'on':5 'rat':12 'sat':4}

# Outline

| Introduction | Realization | **Custom C-functions** | Rating comparison | Conclusion |
|:---|:---|:---|:---|:---|
| oo | oooo | o●o | oooooo | ooo |

# Adding your custom C-functions to PostgreSQL

### Prerequisites

- Developer version of PostgreSQL
- Installation of make
- Root privilege on database

### Folder structure
```
Extension
  ├── function.c
  ├── Makefile
  ├── function.control
  ├── function--1.0.sql
  ├── README.function
```

### Steps

(1) make install

(2) CREATE EXTENSION "extension"

## Example

```
 1  #include "postgres.h"
 2  #include "fmgr.h"
 3  #include "utils/geo_decls.h"
 4
 5  #ifdef PG_MODULE_MAGIC
 6      PG_MODULE_MAGIC;
 7  #endif
 8
 9  PG_FUNCTION_INFO_V1(add_one);
10
11  Datum
12  add_one(PG_FUNCTION_ARGS)
13  {
14      int32   arg = PG_GETARG_INT32(0);
15      PG_RETURN_INT32(arg + 1);
16  }
```

# Outline

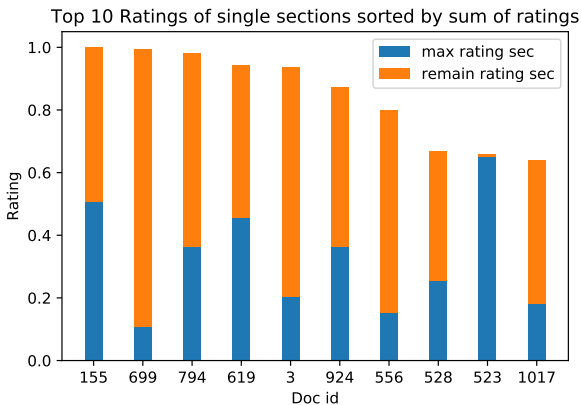| Introduction | Realization | Custom C-functions | Rating comparison | Conclusion |
| :-- | :-- | :-- | :-- | :-- |
| oo | oooo | ooo | o●oooo | ooo |

## Idea

- Originates from a misunderstanding
- Thought the task is to rank whole wiki pages
- User wants the best section and not the "best" document
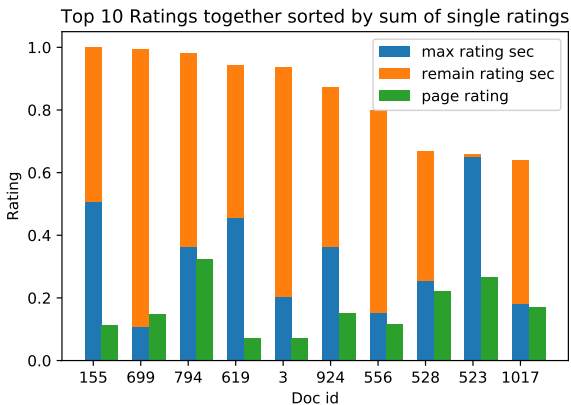- So how is the relationship between page and section ranking

### Calculation of Rating

- **section**: rating / num_words_of_section
- **page**: sum_of_ratings / num_words_of_page

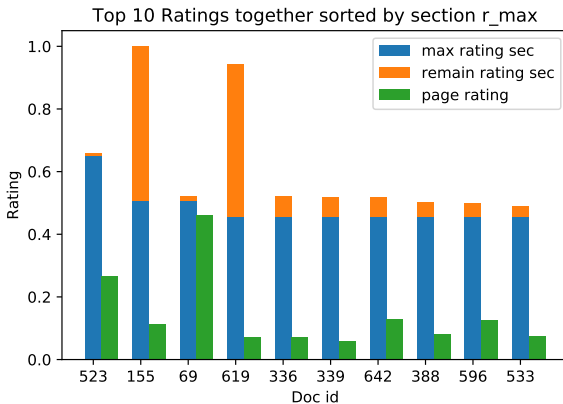# Query:"game", sorted by sum of section rankings



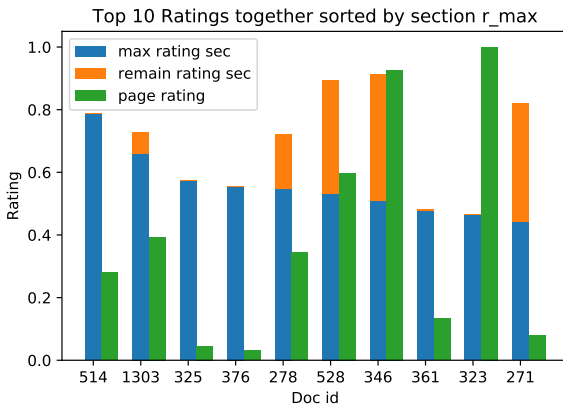Top 10 Ratings of single sections sorted by sum of ratings

# Query:"game", adding the rank for the whole page

# Query:"game", ordered by max section rating



Top 10 Ratings together sorted by section r_max

Introduction
○○

Realization
○○○○

Custom C-functions
○○○

Rating comparison
○○○○○●

Conclusion
○○○

# Query:"game AND team AND ball", ordered by max section rating



Top 10 Ratings together sorted by section r_max

# Outline

1. **Introduction**

2. **Approach and realizations**

3. **Custom C-functions in PostgreSQL**

4. **Rating sections vs. rating pages**

5. **Conclusion**

**Introduction**
oo

**Realization**
oooo

**Custom C-functions**
ooo

**Rating comparison**
oooooo

**Conclusion**
o●o

# Conclusion and future work

### Conclusion

- Ratings for sections and page return total different results
- Tsvector has a lot of potential
- PostgreSQL is easy customizable

### Future work

- Improve the rating algorithm with tf idf information (ts_stat)
- Tests on big datasets

**Introduction**
oo

**Realization**
oooo

**Custom C-functions**
ooo

**Rating comparison**
oooooo

**Conclusion**
oo●

Questions

# **Questions**