

Vision-based Web Page Segmentation

Bc. František Maštera, xmaste02@stud.fit.vutbr.cz

Supervisor: doc. Ing. Radek Burget, Ph.D.

Introduction

WPS (Web Page Segmentation) usage

- Data mining, document indexing
- Assistive technologies

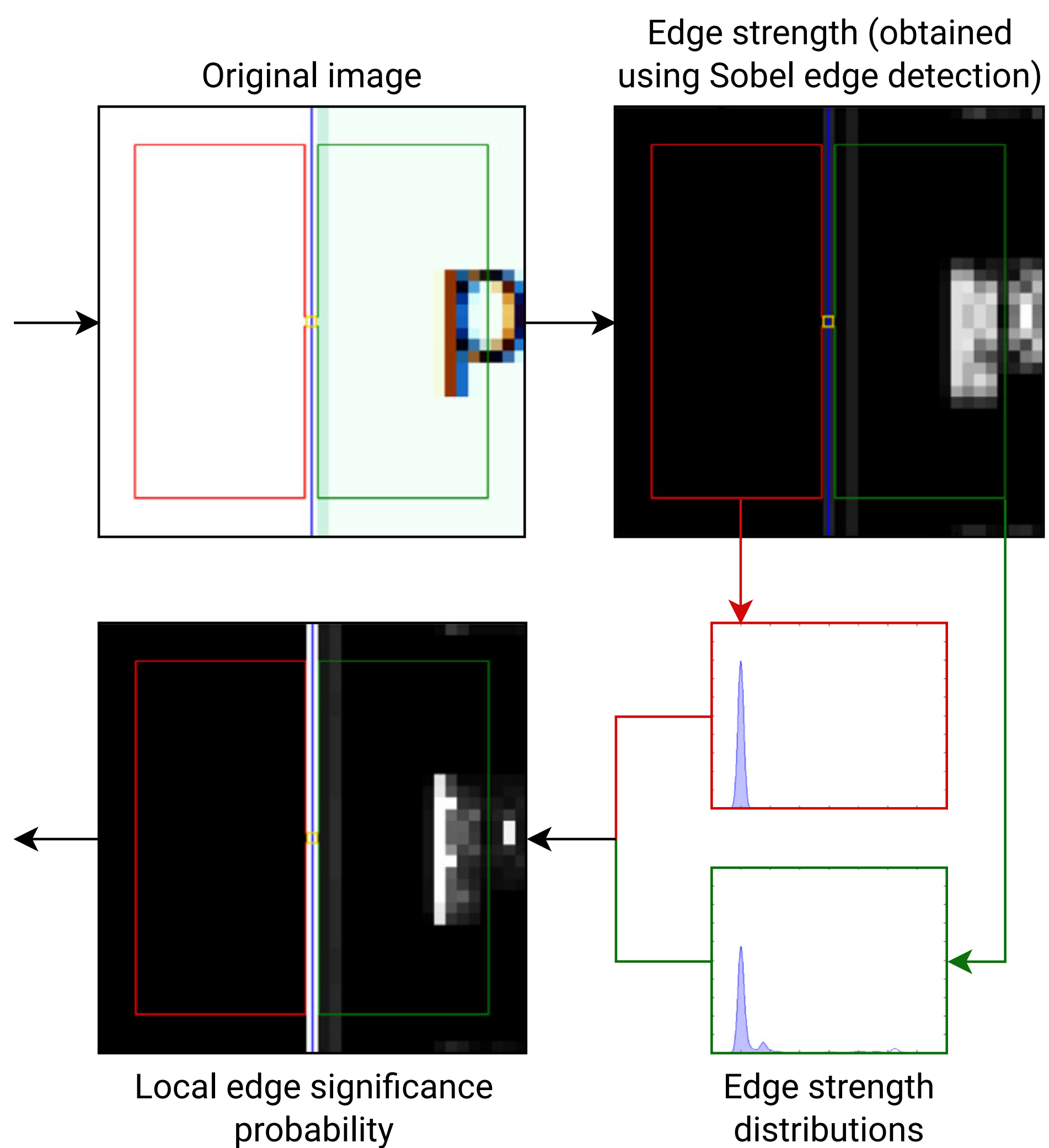
Many WPS algorithms

- With **different approaches**, strengths & weaknesses

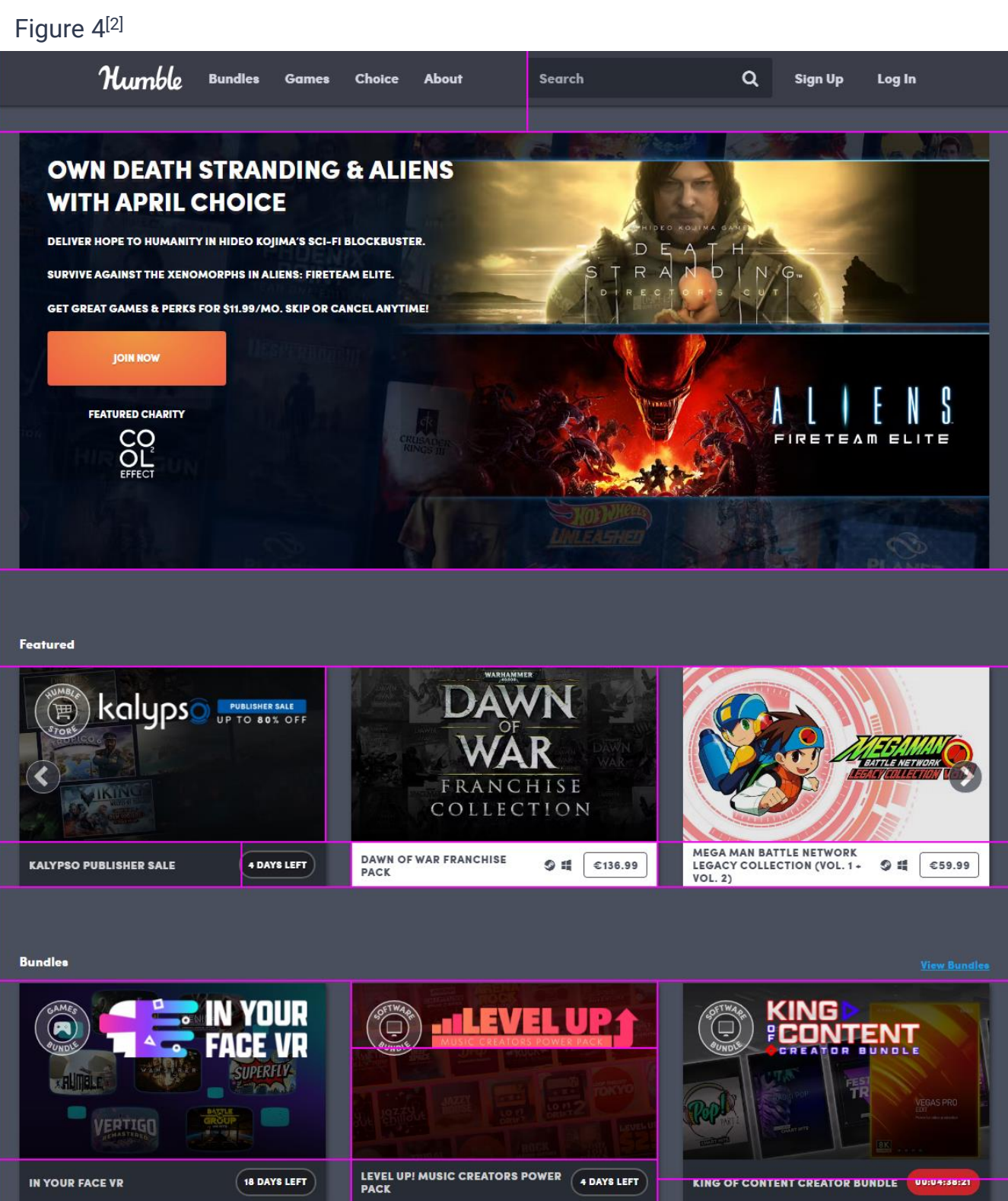
FitLayout framework

- Tools for **WPS evaluation** and further research
- Suite of already implemented WPS algorithms
- Goal: **Extend FitLayout** with another one

Locally Significant Edges



Output Preview



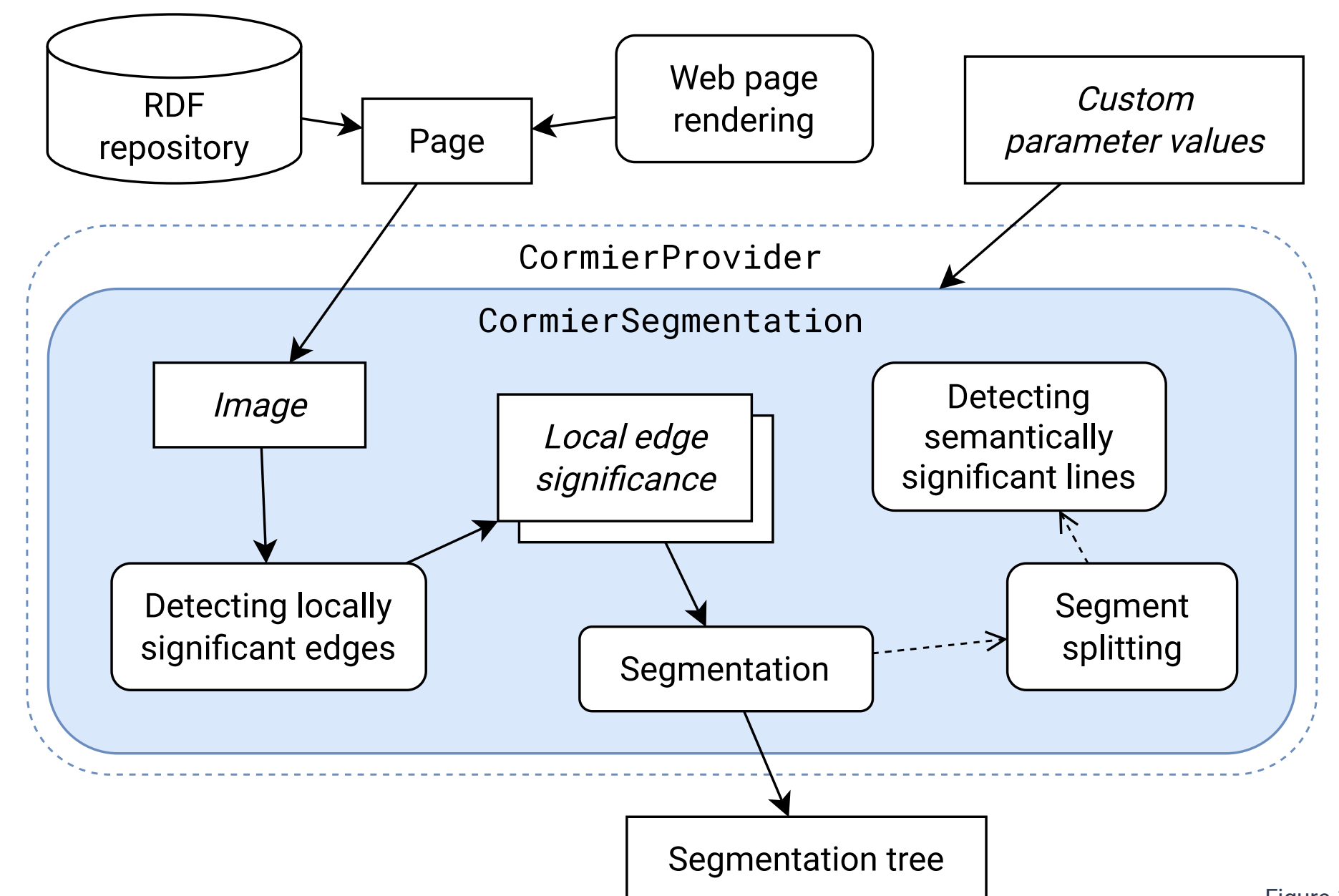
Hierarchical structure

- Main sections:
 - Navigation bar, main section, featured items, ...
- Separate items
- Item sections:
 - Banner, details

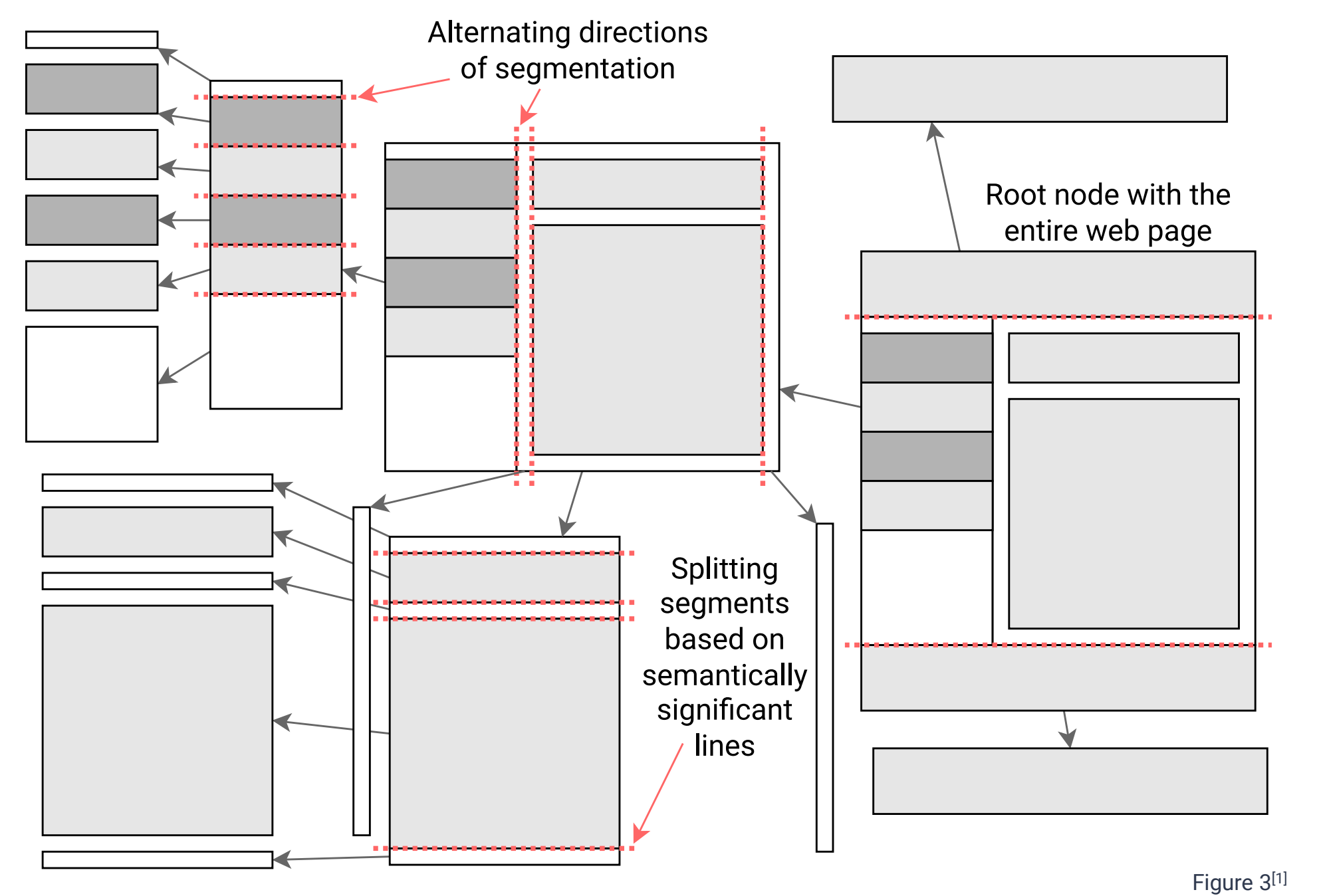
Imperfections

- Inconsistencies
- Poor segmentations

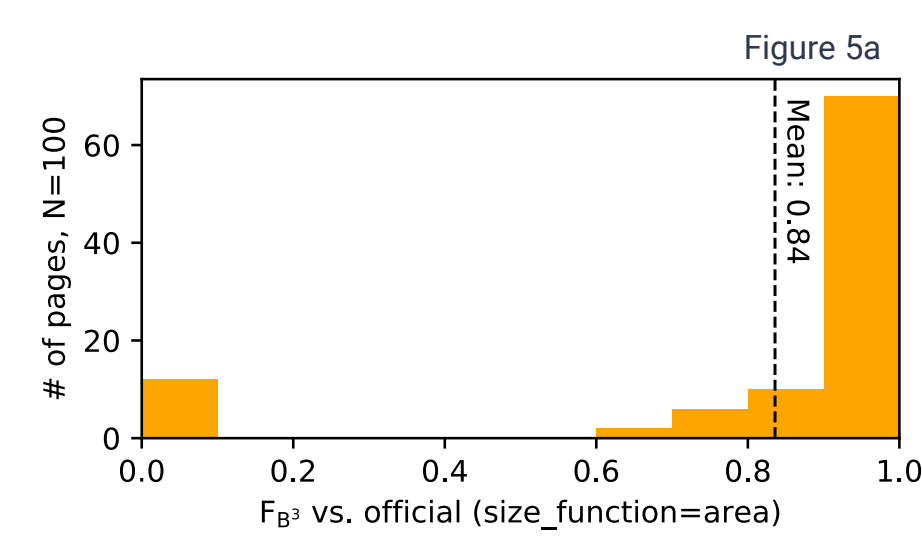
Integration & Pipeline



Final Segmentation

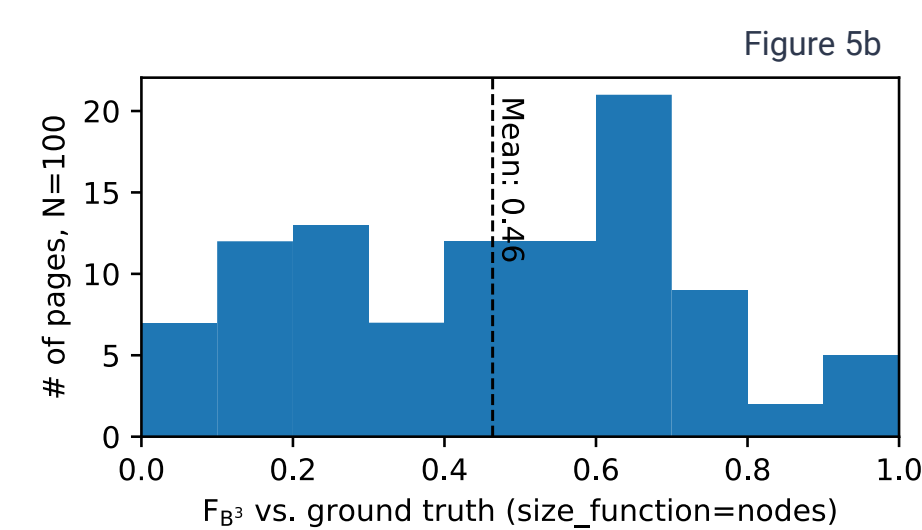


Testing & Evaluation



Testing against public impl.

- Most results meet expectations
- Exceptions with differences



Parameter evaluation

- Affects granularity, speed, ...
- Optimal values for most pages
- Improve the results by **23.13%**

