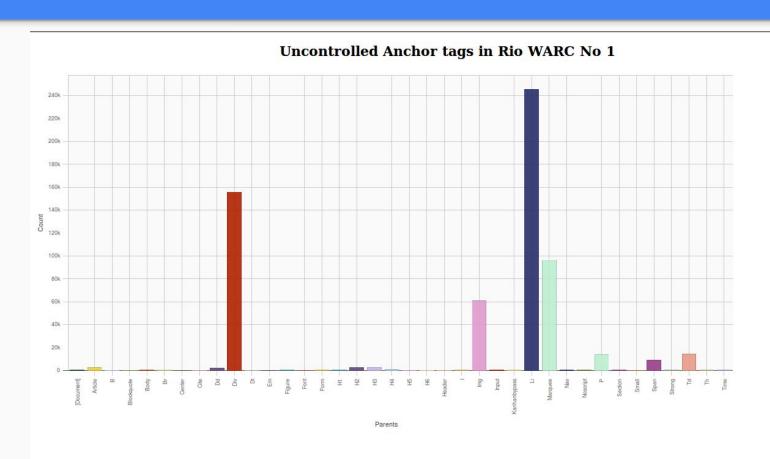# Link Ranking group

Gregory Wiedeman (University at Albany, SUNY)
Mindaugas Vidmantas (The British Library)
Peter Webster (Independent Scholar & Consultant)
Kees Teszelszky (National Library of the Netherlands)
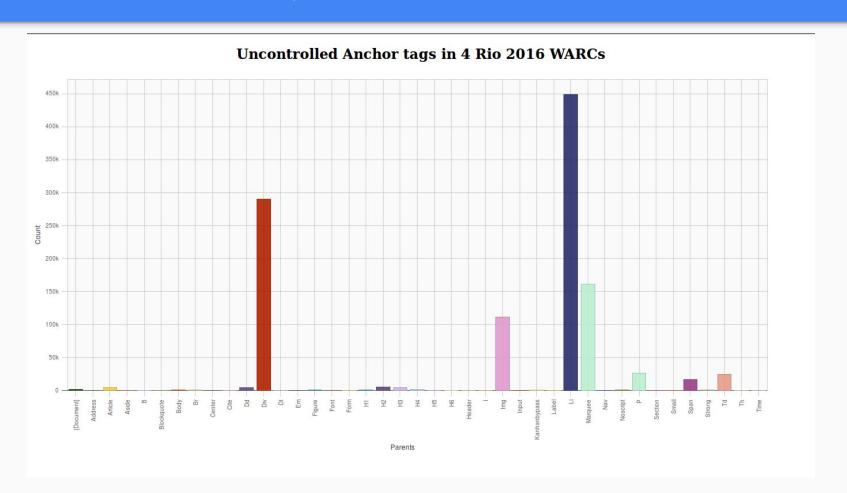Richard Deswarte (University of East Anglia)

# Are All Links Created Equal?

- WarcBase scripts to export manageable raw HTML
- Load into BeautifulSoup Python library
- Look for <a> parents
- Should we weigh links with certain parents more during during link analysis?
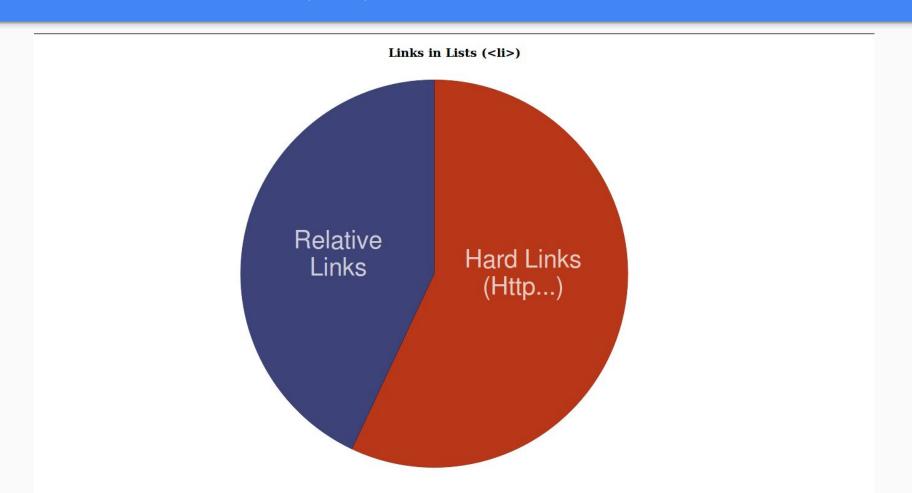- Are navigational links (<li>) different that content links (<p> or <div>)?

Uncontrolled Anchor tags in Rio WARC No 1

# 4 WARC Link Parent Element Type Distribution



Uncontrolled Anchor tags in 4 Rio 2016 WARCs

# Relative vs hardcoded in Rio (2016)



**Links in Lists (<li>)**

Relative Links

Hard Links (Http...)

# Relative vs hardcoded in Rio (2016)



**Links in Content (&lt;div&gt; or &lt;p&gt;)**

Relative Links

Hard Links (Http...)

# CPP Link Parent Element Type Distribution



Uncontrolled Anchor tags in 4 2005 CPP WARCs

# Relative vs hardcoded in CPP (2005)



**Links in Tables (<td>)**

Hard Links (Http...)

Relative Links

# Absolute and Relative Paths in 2005 CPP vs. 2016 Rio WARCs



**Relative or Absolute Links in all 2005 CPP WARCs**

Hard Links (Http...)

Relative Links

**Relative or Absolute Links in all Rio 2016 WARCs**

Relative Links

Hard Links (Http...)