

How big can a static site be?

Staticizing a census database

Martin Holmes and Greg Newton
University of Victoria Humanities Computing and Media Centre



University
of Victoria

Humanities Computing
and Media Centre



Project Endings and static websites: Why?

According to our research the long-term viability of DH projects is tenuous at best

- More than half of all DH projects do not have long-term preservation plans
- More than half of all DH projects do not include planning for an endpoint
- Only 10% of all DH projects consider their documentation to be adequate



The long-term viability of DH projects is tenuous at best



- Over 20% of DH projects stop working due to software obsolescence
- Responsibility for long-term maintenance falls to PI or nobody in nearly half of all DH projects
- Ongoing funding for maintenance is a major obstacle



Project Endings and static websites: Why?

Restoring/updating old projects is time-consuming and costly.

Project Endings offers practical solutions for obviating the need for long-term maintenance.

Is your project viable in the long term?

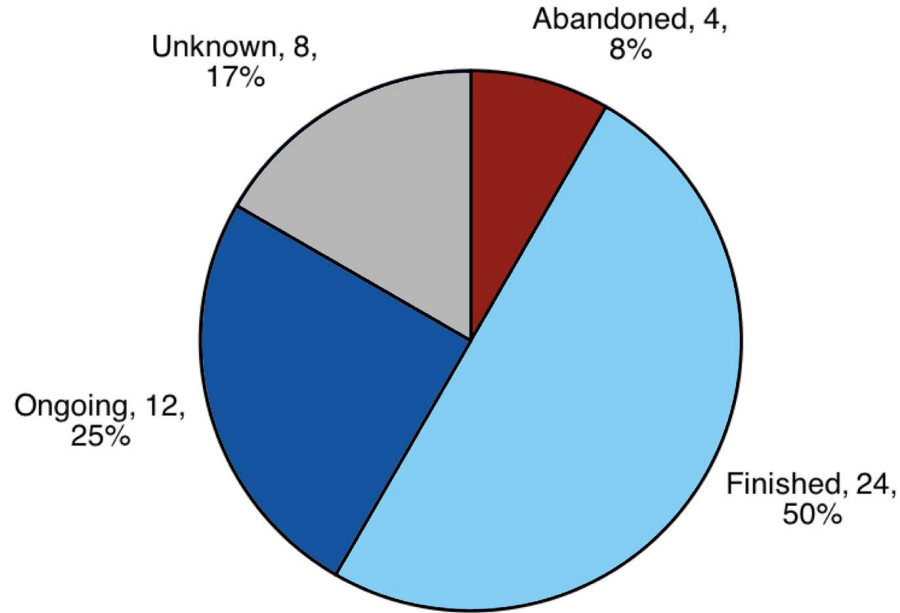
<https://hcmc.uvic.ca/endings/questionnaire.htm>



University
of Victoria

Humanities Computing
and Media Centre

Project status



2019 status of 48 DH2005 projects that had a web component

[The Final Death\(s\) of Digital Scholarship. Davis \(2019\)](#)



University
of Victoria

Humanities Computing
and Media Centre

Project has the most which is most content with the least.

– Diogenes

Be like Diogenes and throw away the cup

- Keep: HTML, CSS, Javascript
- Avoid: external dependencies
- Discard databases
- Say no to CMSs



University
of Victoria

Humanities Computing
and Media Centre

Project Endings and static websites: How?

Perceived realities

- Applications needs maps, graphs and all sorts. Can I avoid third-party tools?
- I need a database! / Where does my data live?
- How do I search my data?
- Can big applications operate with such constraints?



Project Endings and static websites: Working examples

The Map of Early Modern London (13,086 pages: mapoflondon.uvic.ca)

The Colonial Despatches (10,826 pages: bcgenesis.uvic.ca)

Digital Victorian Periodical Poetry (20,685 pages: dvpp.uvic.ca)



Our main focus has been *doing this for real to demonstrate that it is practical.*

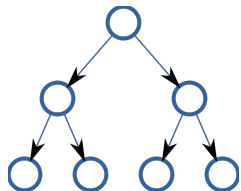


University
of Victoria

Humanities Computing
and Media Centre

But is there a limit?

Document collections such as digital editions need:



browsability (drill down through a hierarchy to find documents)



searchability (look for text, filter results based on dates, document types and so on)

Is there a limit on the scale of project which can be staticized, and if so, what is it?

We searched for a candidate to push the boundaries.



VIHistory (Vancouver Island History)

- PostgreSQL/PHP project
- 15 years old
- Primarily census data from Vancouver Island
- Census records from 1871, 1881, 1891, 1892*, 1901, 1911
- Around 150,000 records
- Associated tables of occupations, familial relationships, locations, addresses, religions, languages, nationalities and more.
- Already partially broken due to PHP and DB updates.



Challenges:

- The quality of the data
- The nature of census data itself
- The organization of the data in the db



Quality of the data

RECEMENT 1881 - CENSUS

RECEMENT 1881 - CENSUS

Province of *British Columbia*

District No. *190*

S District *Johnson St. West*

SCHEDULE No. 1 - Nominal Return of the Living.

TABLEAU No. 1 - Énumération des Vivants.

PAGE *92*

NOMENCLATURE OF THE HOUSEHOLD	NAME	SEX	AGE	Date of Birth	Place of Birth	RELIGION	OCCUPATION	Profession, Occupation or Trade	Marital Status	EDUCATION	Date of Immigration																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																														
1	2	3	4	5	6	7	8	9	10	11	12																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																														
1						<i>Calhoun Ann</i>	<i>F</i>	<i>2</i>	<i>✓</i>	<i>N.B.</i>	<i>Presbyterian</i>	<i>Irish</i>																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																													



Province of *Johnson*
Eager Lane 12

District No.

190

S District

RECENSSEMENT 1901 - CENSUS

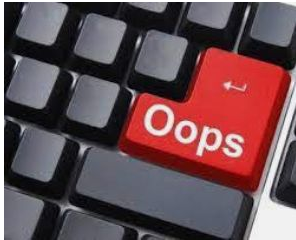
SCHEDULE No. 1 - Nominal Return of the Living.
 TABLEAU No. 1 - Énumération des Vivants.

PAGE 92

Description of the Census of the Census						NAME	SEX	AGE	Date within last census	Country or Province of Birth	RELIGION	ORDRE	Profession, Occupation or Trade	Marital or Single	Instruction				Date of Census and Remarks
French	English	Native	Native	Native	Native										Reading	Writing	Other	Other	
Description of the Census of the Census						NAME	SEX	AGE	Date within last census	Country or Province of Birth	RELIGION	ORDRE	Profession, Occupation or Trade	Marital or Single	Instruction				Date of Census and Remarks
French	English	Native	Native	Native	Native										Reading	Writing	Other	Other	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
						<i>Belcast Ann</i>	<i>F</i>	<i>2</i>		<i>N.B.</i>	<i>Presbyterian</i>	<i>Irish</i>							<i>May 12. Eager Lane 12</i>
						<i>Carlton Horatio</i>	<i>M</i>	<i>53</i>		<i>N.B.</i>	<i>Catholic</i>	<i>German</i>	<i>Y.M.C.A.</i>	<i>M</i>					
						<i>Wagon Ann</i>	<i>F</i>	<i>50</i>		<i>N.B.</i>		<i>Irish</i>							
						<i>Wagon George</i>	<i>M</i>	<i>29</i>				<i>German</i>							<i>Wagon George, age 29.</i>
						<i>Edwin H.</i>	<i>M</i>	<i>21</i>					<i>Shoreman's</i>						
						<i>Horace</i>	<i>F</i>	<i>19</i>					<i>brew maker</i>						
						<i>William</i>	<i>M</i>	<i>17</i>											
						<i>Minnie</i>	<i>F</i>	<i>15</i>					<i>dressmaker</i>						
						<i>Wagon</i>	<i>M</i>	<i>6</i>											
						<i>Robert</i>	<i>M</i>	<i>1</i>		<i>N.B.</i>									
						<i>White James Al.</i>	<i>M</i>	<i>33</i>		<i>Scotland</i>	<i>Presbyterian</i>	<i>Scottish</i>	<i>Seton Kirk</i>	<i>M</i>					
						<i>Minnie</i>	<i>F</i>	<i>21</i>		<i>USA</i>		<i>American</i>							
						<i>Hargrave</i>	<i>M</i>	<i>4</i>		<i>N.C.</i>		<i>Scottish</i>							
						<i>Minnie</i>	<i>F</i>	<i>2</i>											
						<i>Henry Al.</i>	<i>M</i>	<i>2</i>											
						<i>Pottinger George</i>	<i>M</i>	<i>56</i>		<i>Scotland</i>	<i>Presbyterian</i>	<i>Scottish</i>		<i>M</i>					
						<i>Isabella</i>	<i>F</i>	<i>52</i>											
						<i>Thomas</i>	<i>M</i>	<i>20</i>					<i>carpenter</i>						
						<i>David</i>	<i>M</i>	<i>15</i>		<i>N.B.</i>			<i>Store Clerk</i>						
						<i>Burns Frank H.</i>	<i>M</i>	<i>26</i>		<i>N.S.</i>	<i>Presbyterian</i>	<i>Irish</i>	<i>Landscaper</i>	<i>M</i>					
						<i>May</i>	<i>F</i>	<i>24</i>		<i>USA</i>		<i>German</i>							
						<i>Harold H.</i>	<i>M</i>	<i>12</i>		<i>N.C.</i>		<i>Scottish</i>							
						<i>Forbes William</i>	<i>M</i>	<i>24</i>		<i>USA</i>	<i>Ch. Eng.</i>	<i>American</i>	<i>Cabinet maker</i>	<i>M</i>					
						<i>Frederic</i>	<i>F</i>	<i>25</i>		<i>USA</i>		<i>English</i>							
						<i>Ellen E.</i>	<i>F</i>	<i>4</i>		<i>N.C.</i>		<i>American</i>							

Opportunities for error

- Original enumerators make mistakes
- Transcribers working from grainy microfiche make mistakes
- Entry into spreadsheets with no data-constraints adds more
- Ingestion from spreadsheets into database creates yet more



So:

Values for gender
(should be “M” or
“F”)

- 0
- 9
- S
- D
- !

Values for hourly
wage (1911):

- \$40.00
- \$100.00
- \$500.00
- \$3,000.00



Challenges of census data

Each census collects information different from the last, because the preoccupations of society and government change.



The 1871 census is obsessed with race:

There are fields for counts of how many of these are in the household:

- white male
- white female
- chinese male
- chinese female
- colored [sic] male
- colored female
- native male
- native female



Chinese houseboys for the Kenneth McKenzie family at Craigflower and Lakehill



In 1891, there's a sudden interest in living conditions:

- building type (shanty, house, hotel)
- construction (wood, brick, stone, longhouse)
- number of floors
- number of rooms



506 Government Street, Hon. John Robson, M.P., residence, 1891



By 1911, the focus is on work and money:



- occupation
- earnings
- employment, employment state, other employment
- weeks working
- hours working
- weekly wage, hourly wage
- insurance (life and health)

Woodworkers Ltd. Sash and Door Factory, 2843 Douglas Street, c. 1912



University
of Victoria

Humanities Computing
and Media Centre

This variability was reflected in the DB structure...

Every census was in a separate table.

Every census table had different fields.

Linked tables for common values (location, nationality, religion) were chaotic because these values were expressed differently from year to year.

Most fields were “text” because the incoming data was so variable.



...and in the web application interface

Every census had a separate search page with a different interface.

There was a global search that attempted to schmush everything together, but omitted one entire census by accident, or perhaps in despair.

Searching was only possible within specific fields, not across an entire record.



Stage 1: convert to XML

We wrote a custom XML schema.

We have datatypes such as “moneyAmountOrGarbage”:

```
<dataSpec ident="vih.moneyAmountOrGarbage" module="vihistory">
  <desc>An amount of money in dollars and optional cents, or some garbage
from the source.</desc>
  <content>
    <dataRef name="token" restriction="(\d+(\.\d\d)?)|(QUERY: .+)" />
  </content>
</dataSpec>
```



One content model to rule them all

All entries use the same content model, with optional components and flexible structures smoothing out the differences between census years:



```

<entry date="1881" legacyId="56137" xml:id="cr_1881_56137">
  <title>Arme, (30), 1881, Victoria (190), Victoria City Johnson Street Ward (B)</title>
  <person gender="Female" ageYears="30" attendingSchool="false">
    <persName>
      <familyName>Arme</familyName>, <givenName/>
    </persName>
    <maritalStatus legacyId="1">Single</maritalStatus>
    <family ref="fam:f_1881_47_147">British Columbia, Victoria, Victoria City Johnson Street
      Ward</family>
    <religion legacyId="980">
      <desc genToken="rel_noneno_religion">None/No Religion</desc>
    </religion>
    <infirmities legacyId="0">Blank or None</infirmities>
    <event type="ownBirth">
      <place legacyId="15100">British Columbia, Canada (Native Indian/First Nations,
        Canada)</place>
    </event>
  </person>
  <location legacyId="47" srcType="fromCensus">
    <censusDistrict legacyId="190" date="1881">Victoria (190)</censusDistrict>
    <censusSubdistrict legacyId="B">Victoria City Johnson Street Ward (B)</censusSubdistrict>
  </location>
  <housing legacyId="2">
    <desc>House</desc>
  </housing>
  <work>
    <occupation legacyId="1057" code="X2100">
      <desc genToken="occ_none_or_unknown" srcType="fromCensus">None or Unknown</desc>
    </occupation>
  </work>
</entry>

```



Payoffs from data conversion

Removed over 500 completely empty records.

Removed about 4,500 mysterious unsourced records, possibly consisting of generated test data.

Normalized and corrected errors in hundreds of records.



Stage 2: Render to HTML

Granularity: 1 census record = 1 page

Browsability/drill-down: 1 data-point (a religion, a nationality, a family) = 1 listing page

Searchability: each census record has many <meta> tags from which staticSearch creates search filters (date, religion, nationality, location, race...)



Result

138,744 individual census record pages

166,528 pages across the site (because of generated listings pages)

Site size on disk: 3.6GB (of which 2.0GB is search index files)

Site build time: 1 hour 46 minutes (on Jenkins CI server)

staticSearch works fine – feel free to try it at the temporary site location:

<https://hcmc.uvic.ca/project/vicensus/>



Conclusions

~~We don't need no stinking database.~~ Even projects that look like a “natural” fit for a database may work better as static sites.

150,000 pages is not a big number for a static site, especially if they're small.

Sophisticated data constraints are easier in XML, so bad data is easier to control.

HOWEVER...



Conclusions (2)

We may be approaching a limit with the current staticSearch, because the index file containing document titles is approaching 16MB. This file has to be downloaded to the client before any search results can be shown, so the first search may appear to take a few seconds longer than subsequent searches.

We can imagine workarounds for this, involving changing the granularity (i.e. aggregating records into larger pages) at the project level, or splitting out the title file at the staticSearch build level.



Acknowledgements & links

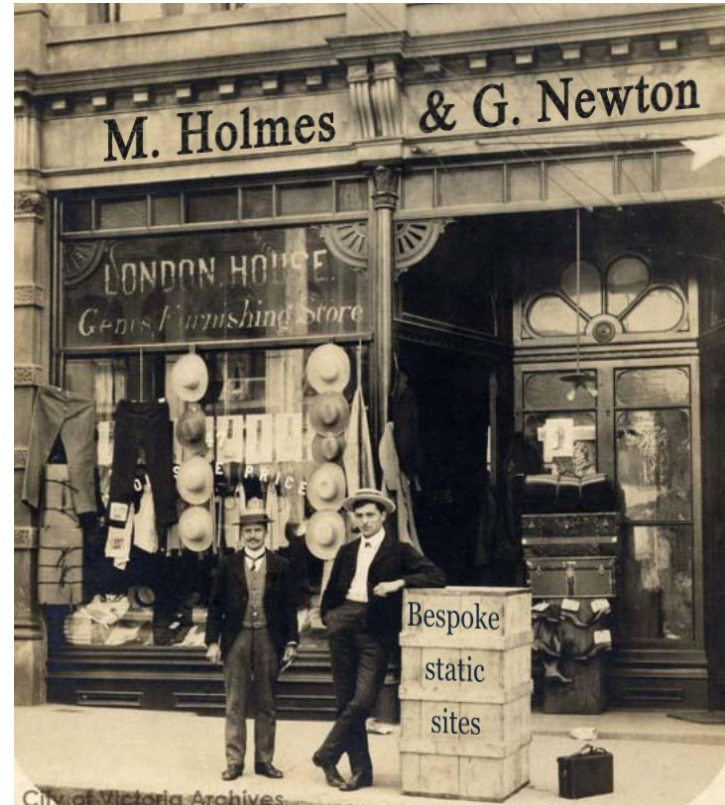
Thank you SSHRC for Endings funding.

Thank you Pat Dunae, VIHistory PI, for clear explanations and infinite patience.

Thank you Joey Takeda for co-authoring staticSearch.

endings.uvic.ca

github.com/projectEndings/staticSearch



University
of Victoria

Humanities Computing
and Media Centre