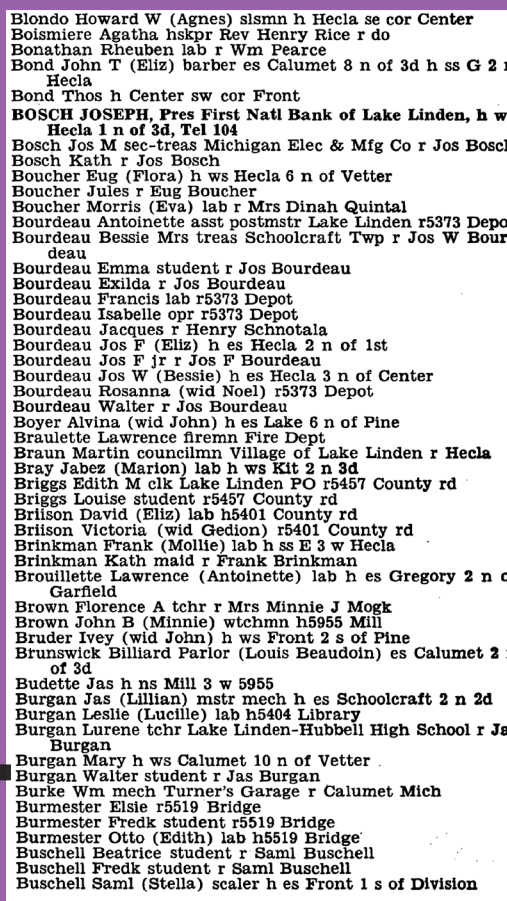


Gary Spikberg, Department of Computer Science  
Ankitha Pille, Department of Computer Science  
Elijah Pass, Department of Computer Science  
Robert Pastel, Department of Computer Science  
Don Lafreniere, Department of Social Sciences

The aim of this project was to create a semi-automated approach for converting text based city directories into geographically referenced points in a GIS. This project contributes to the larger Keweenaw Time Traveler project which is creating a 3D spatial model of the industrial mining communities in Michigan's Copper Country complete with a digital reconstruction of the built and social environments from 1880-1950



After scanning the directories we use ABBYY Finereader Pro to crop to the names, and read them using the Optical Character Recognition (OCR) function. From here we save a text only copy of the directory, as well as the cropped image. Both will be used later.

To allow for rapid importation and searchability of the city directory data, we needed to parse information about each resident from the directory into individual segments.

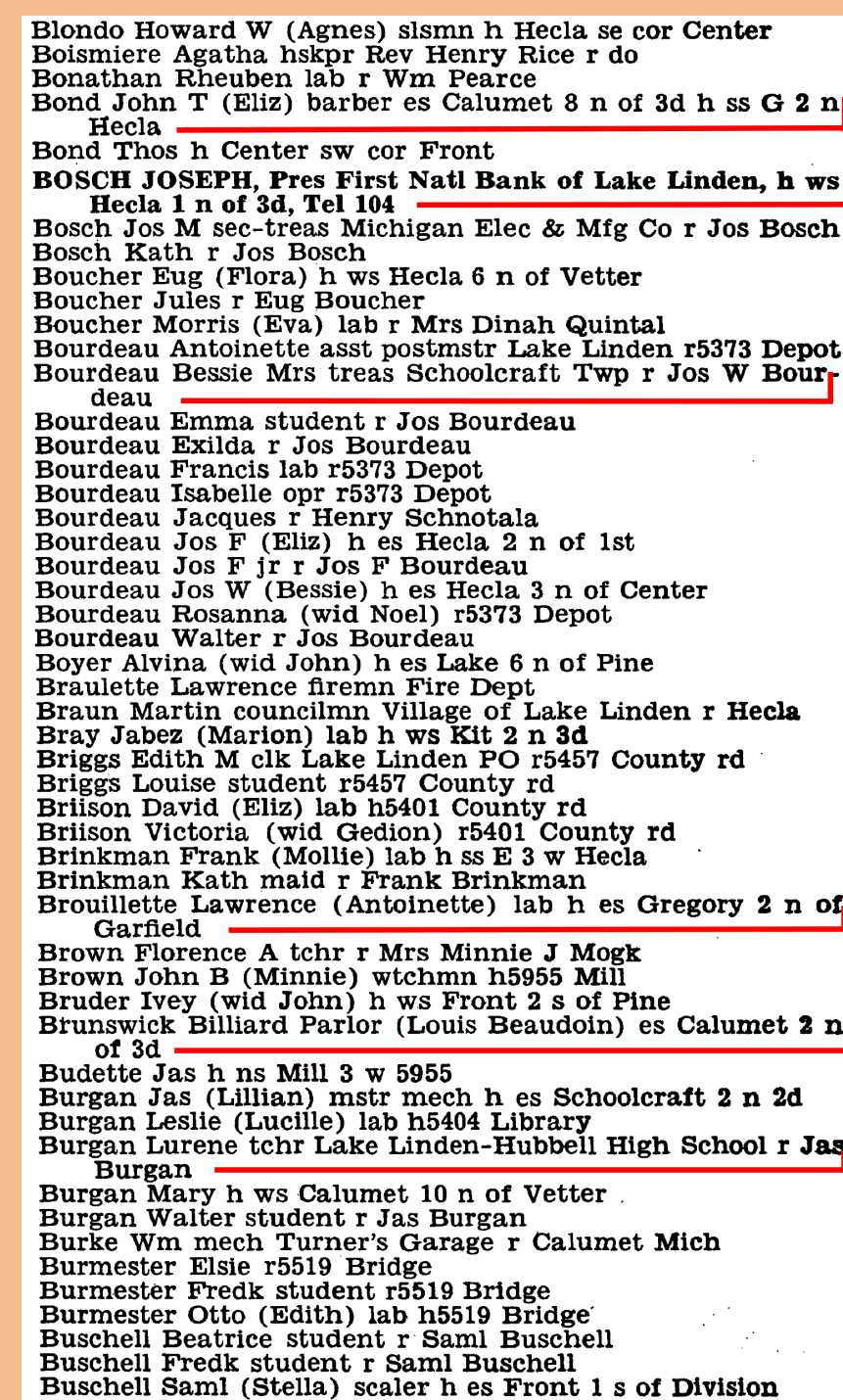
Occupations, workplaces, residential addresses, and housing tenure are captured as well. The blue underlined areas noted at left are occupations.

In addition to exact civic addresses, there are also relative locations for communities that had not yet received home mail delivery . Those were captured by directions (n, e, s, or w) as well as by the name of who they lived

[illegible]

After parsing, we geocode the parsed addresses with a previously completed GIS feature class created by digitizing all of the buildings noted on Sanborn Fire Insurance Plans. The feature class includes the civic addresses and is updated every 10 years roughly coinciding with the decennial census. Several passes through the data are required adjusting the script to account for different tolerances for spellings of streets and street name changes.

To clean the text file before it can be parsed all of the multiline entries need to be placed on one line. A manual cleaning of OCR errors must be complete. For example, 'W's tend to be recognized as 'VV' or 'V' due to faded print. This is especially difficult when working from microfilm copies of city directories.



We created a Groovy script, using regular expressions, to parse the entries automatically. We first isolated words starting with a capital letter (to capture proper names) and recorded middle or first initials where possible.

This process was evolved by and tailored to each different set of directories, mostly due to the entry formats changing in both different cities, and different years. The directory to the right contains punctuation after the name, occupation, address. This made it easier to isolate each segment.

A small output sample:

outUnitNum	outFirstName	outLastName	outMId	outRfMs	outWId	outSpouse	outJob	outWorkPlc	outC
2	Aashia	MaryAnn	null	Mrs	wed	jos	null	null	r
3	Alan	Victor	null	Mr	wed	jos	lab	h	h
4	Alvin	Leo	null	null	null	lab	lab	h	h
5	Alvin	Victor	null	null	null	ida	null	h	h
7	Amos	Nancy	null	Mr	wed	wed	lab	h	h
8	Andler	Harry	J	null	null	tchr	Lab	Lake Linden	r
9	Archie	Barnette	null	Mr	Wed	lab	lab	h	h
13	Archambeau	Ida	null	null	null	null	null	h	h
12	Asselin	Albert	J	null	null	Lake Linden	Unl	h	h
13	Atkins	Normand	null	Mr	Null	Null	Null	h	h
14	Aubin	Jos	H	null	Null	Mary	lab	h	h
15	Aubin	Lawrence	null	null	Null	Pauline	lab	h	h
16	Aubin	Nelson	null	Mr	Null	Rebecca	lab	h	h
19	Ausette	Raymond	null	null	null	lab	bkr	First Natl Bar	r
20	Baldwin	Alphonse	null	Mr	Null	lab	lab	h	h
23	Baldman	Amanda	null	null	null	etk	etk	A P Tea C	r
25	Ballgaron	Jos	null	Mr	lab	lab	lab	h	h
26	Barnes	Alva	null	Mr	Null	pres	lab	Age of Lake	h
28	Bart	Alphonse	null	null	Null	Genevieve	lab	h	h
30	Bartl	Arth	null	Mr	Null	Mary L	stumn	Nelson Canth	h
31	Bart	Henry	Henry	Mr	Null	Null	Null	h	h
34	Bartl	Jas	null	null	Null	councilmn	Village of Lak	h	h
35	Bartl	Jos	null	null	Null	Victoria	brkm	unl	h
39	Bartl	Mose	null	Mr	Null	Victoria	lab	h	h
40	Bartl	Oliver	null	Mr	Null	Alina	justice of the Shoolcraft	h	h
41	Barnette	Pubert	null	Mr	wed	Null	Null	h	h
43	Bartl	Walter	null	Mr	Null	Delia	Null	h	h
45	Barnette	Urie	null	null	Null	Mary	lab	unl	h
46	Barden	John	null	Mr	Null	Etia	lab	h	h
47	Beauchaine	Jos	null	Mr	Null	Clara	lab	h	h

The final step requires researchers to review matches to ensure accuracy and to manually place relative locations on the map.

The final step requires researchers to review matches to ensure accuracy and to manually place relative locations on the map.

City and Year	# of Directory Entries	Total Parsed	% Parsed	Auto Geocoded	Manual Geocoded	% Geocoded
Hancock 1939	4504	4138	91.87%	2494	1095	86.73%
Hancock 1930	3788	3450	91.08%	2069	288	68.14%
Laurium 1930	2575	2459	95.49%	2162	In Progress	87.92%
Laurium 1916	4529	4259	94.04%	1827	In Progress	42.90%
Houghton 1930	3397	3116	91.73%	1457	In Progress	46.76%
Houghton 1916	3606	3142	87.12%	889	In Progress	28.29%