Retrieving a Sorted List of Hundreds of Closest Relatives from FamilySearch Family Tree in Seconds

Family History Technology Workshop Brigham Young University March 20, 2014

Ben Baker bakerb@familysearch.org



© 2013 by Intellectual Reserve, Inc. All rights reserved.



Weighted Relationship Distance

$$\mathsf{WRD}_{(\mathsf{g}, \mathsf{c}, \mathsf{m})} = \alpha(|\mathsf{g}| + 1) \, \mathrm{e}^{\beta \mathrm{c}} \, \mathrm{e}^{\gamma \mathrm{m}}$$

- g Generational or "vertical" distance Number of generations from base person
- c Collateral or "horizontal" distance
 - Minimum generations to a closest common ancestor
- m Marriage distance
 - Number of marriages between base person
- α , β , γ Weighting factors to control growth rates

2013 Family History Technology Workshop – Beyond the Relationship Calculator Using a Weighted Relationship Distance Metric to Prioritize, Categorize and Visualize Relatives http://fht.byu.edu/prev_workshops/workshop13/papers/baker-beyond-fhtw2013.pdf





Same Sample Relatives in Table Format in Sorted WRD Order



Relative Description	g	С	m	WRD
Base Person	0	0	0	1.0
Father/Mother	1	0	0	2.0
Brother	0	1	0	2.72
Grandfather/Grandmother	2	0	0	3.0
Son	-1	0	0	3.52
Wife	0	0	1	4.14
Aunt/Uncle (sibling of parent)	1	1	0	5.44
Husband of Grandmother	2	0	0.5	6.10
1 st cousin	0	2	0	7.39
Father/mother-in-law	1	0	1	8.27
Niece (daughter of brother)	-1	1	0	9.57
Brother/sister-in-law (spouse of sibling or sibling of spouse)	0	1	1	11.25
Daughter-in-law	-1	0	1	14.56
Aunt/Uncle (spouse of sibling of parent)	1	1	1	22.49
Niece (daughter of spouse's brother)	-1	1	1	39.59
Sister-in-law (wife of spouse's brother)	0	1	2	46.53



WRD(-1,18,2) = 3,955,848,641 Wife of 4th cousin 7 times removed of wife of 6th cousin 6 times removed

FamilySearch

Relationship Calculator Deficiencies

- 1. A relationship calculation must be initiated between two persons in an ad hoc manner and repeated for a different set of two persons
- 2. It is not possible to sort a list of arbitrarily related persons by closeness
- 3. Relationship calculations through marriages are not possible
- 4. Data to perform the relationship calculations must be pre-calculated and stored to be performant enough for on-demand calculation in an enterprise system



Apache Cassandra-Based FamilySearch Family Tree





Proposed System

RESTful web service endpoint added to 5-node Cassandra based test system

Method / Description

GET relatives/{id}

Retrieve the closest relatives to a person up to a specified maximum weighted relationship distance.

Parameters

double maxDistance – Optional query parameter (default 10.0) to specify the maximum distance to return relatives of the person up to.

Returns

A list of RelativePerson objects, sorted by increasing distance

Future methods planned to retrieve closest common ancestor and closest relation between two people.



RelativePerson Object

```
"name": "George Dean Cockle",
"id": "KWJX-VMC",
"lifespan": "1901 - 1970",
"fatherIds": [
  "K2V7-PV5"
Ι,
"motherIds": [
  "K2V7-P92"
1,
"spouselds": [
  "MM8P-MGS",
 "KWBB-7G9"
1,
"familyName": "Cockle",
"givenName": "George Dean",
"gender": "MALE",
"childlds": [
],
```

```
"relative": {
 "id": "K2V7-PV5",
 "name": "John Cockle",
 "lifespan": "1852 - 1921",
 "relativeRole": "CHILD"
"weightedRelationshipDistance": {
 "starRanking": 8,
 "simpleRelationshipDistance": "3",
 "weightedRelationshipDistance": "8.155",
 "generationDistance": 2,
 "collateralDistance": 1,
  "marriageDistance": "0"
},
"relativeDescription": "Great Uncle"
```



Performance Results

Maximum WRD	Num Persons	Mean Execution Time
5.0	43	1.95s
10.0	167	22.1s
15.0	554	89.8s

Performance tuning is still necessary to make this service more useful for production use.

For comparison purposes, the average time to retrieve data for a pedigree containing 10-15 persons in FamilySearch Family Tree is typically about 3-5 seconds. The Eureka team has show subsecond pedigree load time involving the same number of persons and loading up to 12 generations in several seconds.





Live Demo



Future Applications Enabled

- Identifying the closest relatives where historical record hints have been identified but not attached yet as sources.
- Promoting e-mail campaigns to point out what others have added to your relatives such as new photos, sources, stories, etc. to draw users back to the site.
- Easily identifying end of line relatives and likely places for successful descendancy research.
- Facilitate LDS temple work for closest relatives first and sharing more distantly related people with others.
- Producing to-do lists sorted by closeness of relation on any task a user may want to undertake (Ex. fixing data anomalies, merging possible duplicates, providing missing data, etc.)
- Automatic watching of relatives via e-mail alerts based on closeness.
- Applications across a set of users (Ex. closeness of relation to the user for all LDS temple submissions or photo uploads on FamilySearch)
- Sorting items such as memories, watched persons, LDS temple reservation lists, etc. in order of those closest to the user.
- Pointing out relatives who have been identified as prominent or have participated in significant events.
- More . . .



Next Steps

- Full transition from Oracle to Cassandra expected to take a while, but expect one-way synchronization sooner
- Enables read-only operations on Family Tree data, including this work
- Intend to petition to utilize this work in production before transition is complete
- May implement a smaller-scale solution in the current system to prove value of WRD metric over current scope of interest service





Additional Q & A

