

Enhancing and validating GeoNames with Digital Nautical Charts: A case study in the mapping of freeform map labels

Dakotah D. Maguire, Jason C. Kaufman,
Alexandre Sorokine, Robert Stewart

Geospatial Science and Human Security
Division

Background

- Temporal and spatial variation in locations vary among multiple sources
 - Cartographic offset vs. precise point locations
 - Name variations and aliases
- Do multiple sources validate each other and add value when combined?

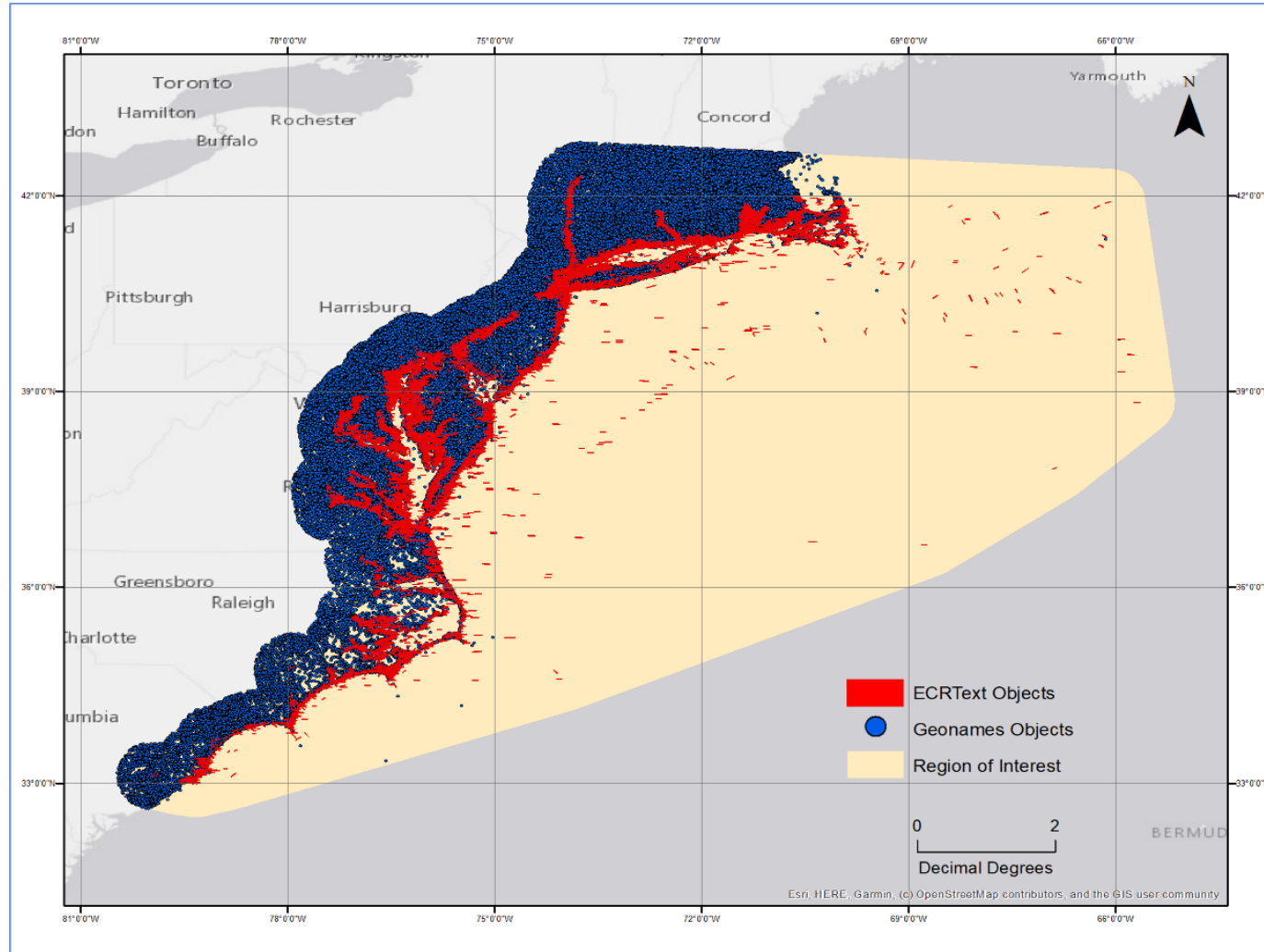
Study question

- Can Digital Nautical Charts (DNC) data enhance GeoNames data by supporting validation, expanding aliases, and filling in missing data?
 - Use textual similarity and spatial proximity
 - Detect instances of ECRTtext within GeoNames

Data

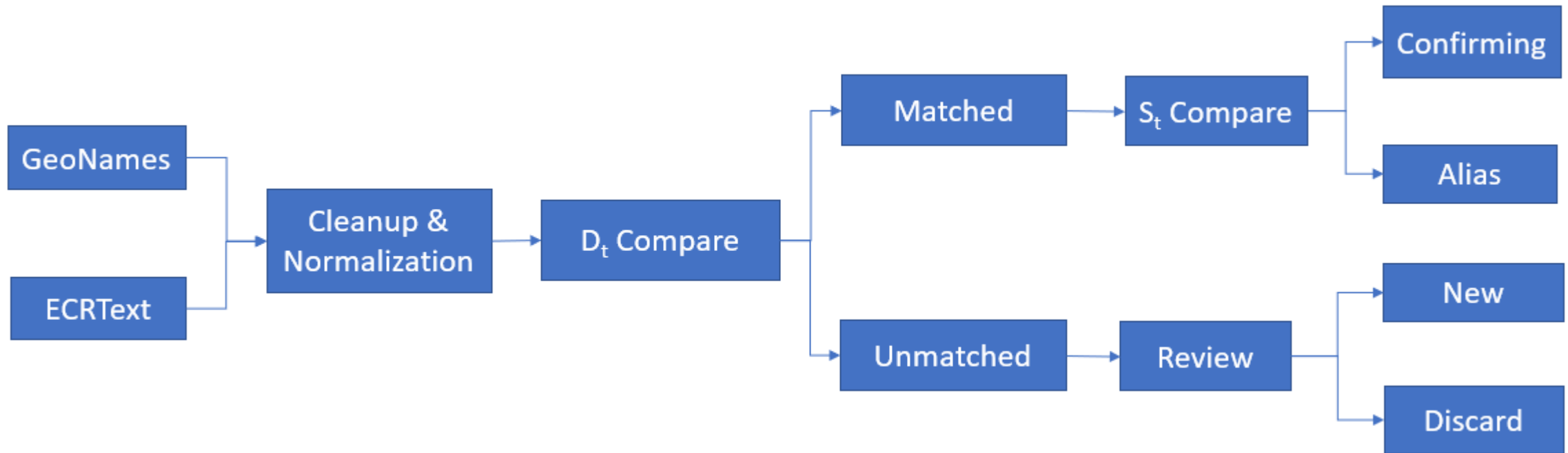
- Digital Nautical Charts is a substantial, long-lived, worldwide vector chart database for ship navigation
 - Developed and maintained by the National Geospatial-Intelligence Agency and the National Oceanic and Atmospheric Administration (NOAA, 2022)
- GeoNames is an open-source gazetteer with over 27 million geographic names

Study area and data



Our region of interest is DNC Region 17 (DNC17), containing 14,224 ECRText objects and 222,772 GeoNames locations along US East Coast from 42° N to 33° N latitude

Methods



Methods: Pairs and similarity

- We can compute the distance between any ECRTText–GeoNames data pair as the closest distance between the GeoNames point and the ECRTText bounding polygon
 - We define a pair as “sufficiently close” if their distance is less than a distance threshold dt
- We use a computationally effective trigrams method to calculate the similarity of the names; results in similarity values in the range of 0–1 (Dunn, 2020)
 - We define a pair as “sufficiently similar” if the similarity value is greater than a threshold st

Methods: Pairs and similarity (continued)

- Semiautomated workflow developed to classify ECRTText
- ECRTText first cleaned and normalized using US Chart No. 1 (NGA, 2019)
- ECRTText and GeoNames within dt distance (a buffered concave hull around ECRTText features) are considered matched pairs and unmatched otherwise
 - Distance threshold value of 0.5° was found to provide reasonable classification outcomes

Methods: Pairs and similarity (continued)

- Matches then compared for text similarities s and classified as confirming or alias
 - Similarity threshold value = 0.8
 - Unmatched pairs manually reviewed to discard results not obviously geographic names; remaining are candidates for new locations

Potential outcomes

- **Confirming:** Here, an ECRTText is part of a sufficiently close pair (distance $\leq dt$) with an exact text–name match (similarity = 1). The interpretation here is that ECRTText very close to a GeoNames location with exactly the same name provides some confirmation that the GeoNames is likely accurate.
- **Alias:** Here, an ECRTText is part of a sufficiently close pair (distance $\leq dt$) with a sufficiently similar name ($st \leq \text{similarity} < 1$). The interpretation here is that ECRTText found very close to a GeoNames location with nearly the same name may well be a new alias (i.e., alternate name) for the already existing GeoNames location.
- **New:** Here, an ECRTText is either too distant from any GeoNames (distance $> dt$) or is part of a pair with significantly different text–name matches (similarity $< st$). The interpretation here is that the ECRTText object is too distant or too different to be reasonably associated with an existing GeoNames location. These may well be candidates for new locations to enhance the GeoNames dataset.
- **Discard:** Here, an ECRTText is too distant from a GeoNames location with a label that is not associated with a geographic name (e.g., Unexploded Ordinance). The interpretation here is these place names are highly unlikely to have any clear benefit to GeoNames.

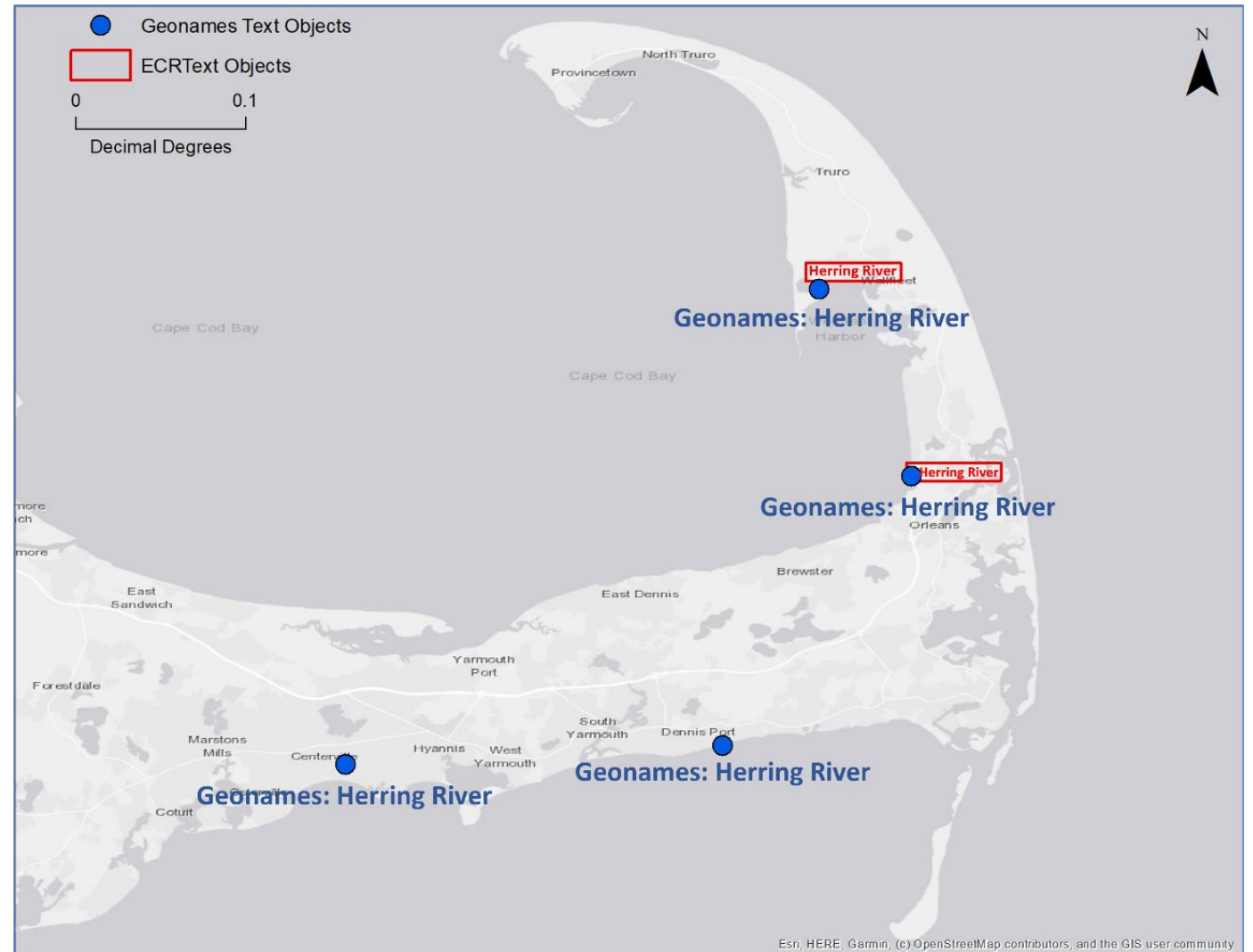
Initial results

- Unique text objects in the study area:
 - 14,224 ECRTText objects
 - 208,203 GeoNames
- Most significant benefit is validation with 87% of ECRTText data classified as confirming
- ECRTText in the Alias (2%) and New (8%) categories has the potential to enhance and extend existing locations

Classification	Total ECRTText Objects
Confirming	12,439 (87%)
Alias	275 (2%)
New	1,103 (8%)
Discard	407 (3%)

Challenges

- Handling crowded places where names are very similar or identical
 - Near things are more related than distant things



Ongoing work

- Conflationary step
 - Duplication across libraries and across scales
- Additional text similarity methods
 - Levenshtein vs. Trigrams performance
- OpenStreetMap as a third data source
 - Potential for more frequent name changes (i.e., added aliases)
 - Duplicate entries, or closely related objects with nearly identical names

Copyright and disclaimer

- This manuscript has been authored by UT-Battelle, LLC, under contract DE-AC05- 00OR22725 with the US Department of Energy (DOE). The US government retains and the publisher, by accepting the article for publication, acknowledges that the US government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this manuscript, or allow others to do so, for US government purposes. DOE will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan (<http://energy.gov/downloads/doepublicaccess-plan>).