

Digitizing EPICC Data: Trials and tribulations in translating 100 year old data

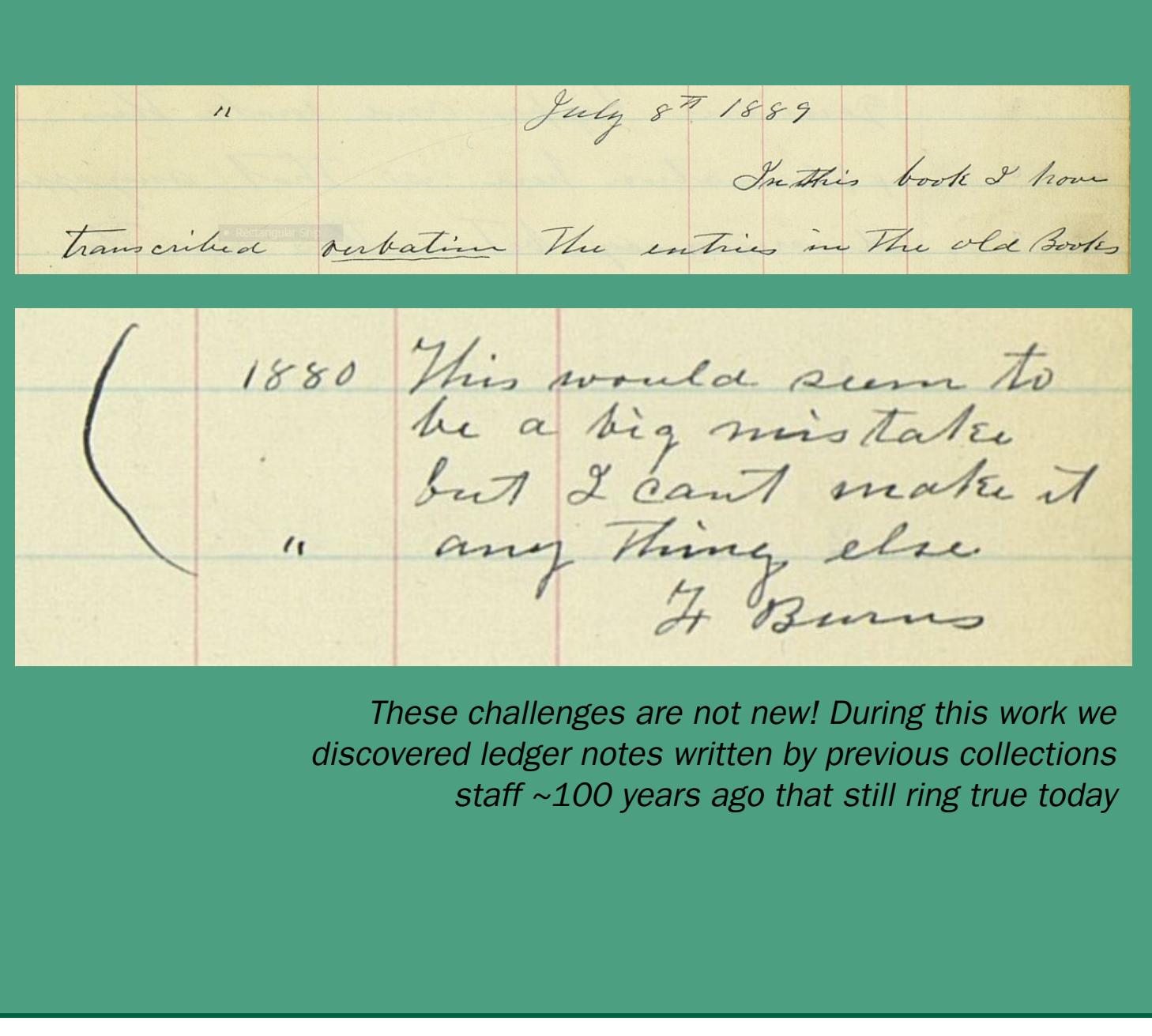
Holly Little, Anna Leary, Alexandra Cano, Adam T. Mansur

The Challenge

The Smithsonian National Museum of Natural History (NMNH) Department of Paleobiology recently completed the first segment of a mass digitization project in support of the Eastern Pacific Invertebrate Communities of the Cenozoic (EPICC) thematic collections network. In collaboration with the Smithsonian Institution Digitization Project Office (DPO), imaging and label transcription was completed for a portion of the Cenozoic Mollusca Collection. Once the labels were transcribed further processing was required to clean and enhance that specimen data. We sought to ensure high quality data for this project through:

1. the development of clear guidelines for documentation and treatment of specific data points;
2. updating records to match current taxonomic, lithostratigraphic, and chronostratigraphic information; and
3. create iterative workflows to maintain extensibility and to capture uncertainty in the data.

A significant challenge for any large collections digitization project is transcribing and cleaning analog information from specimen labels. The following considerations, methods, and tools enabled us to face these challenges of translating analog collections data of over a hundred years old into high quality records following modern standards for biodiversity information.

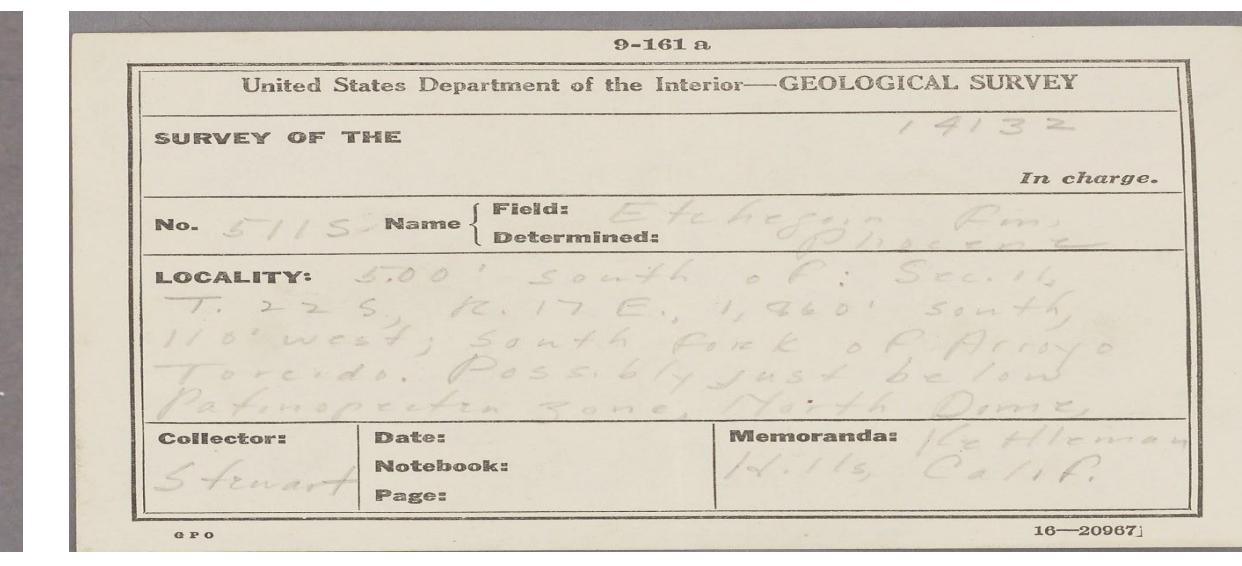
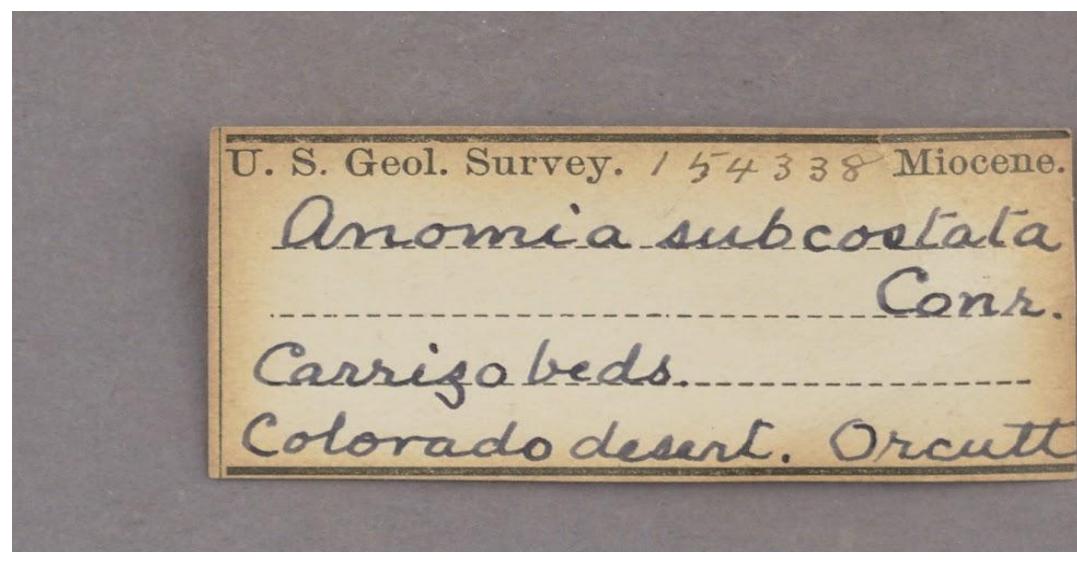
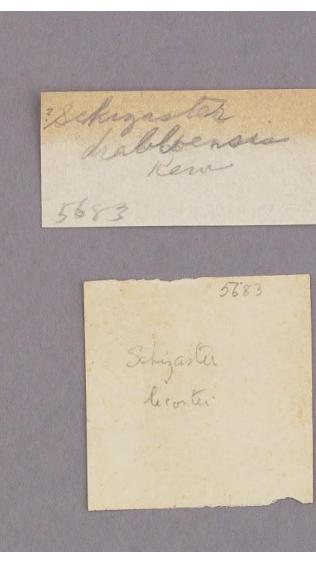
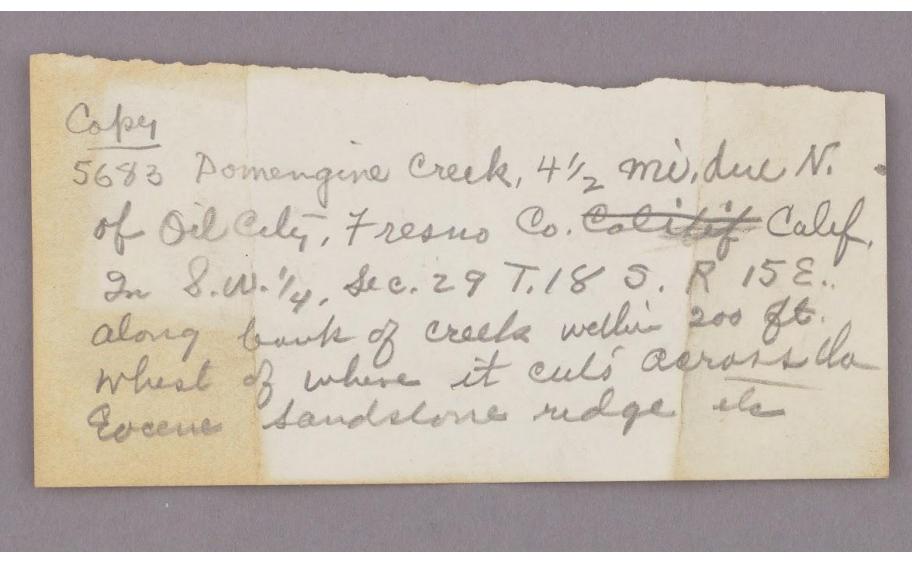
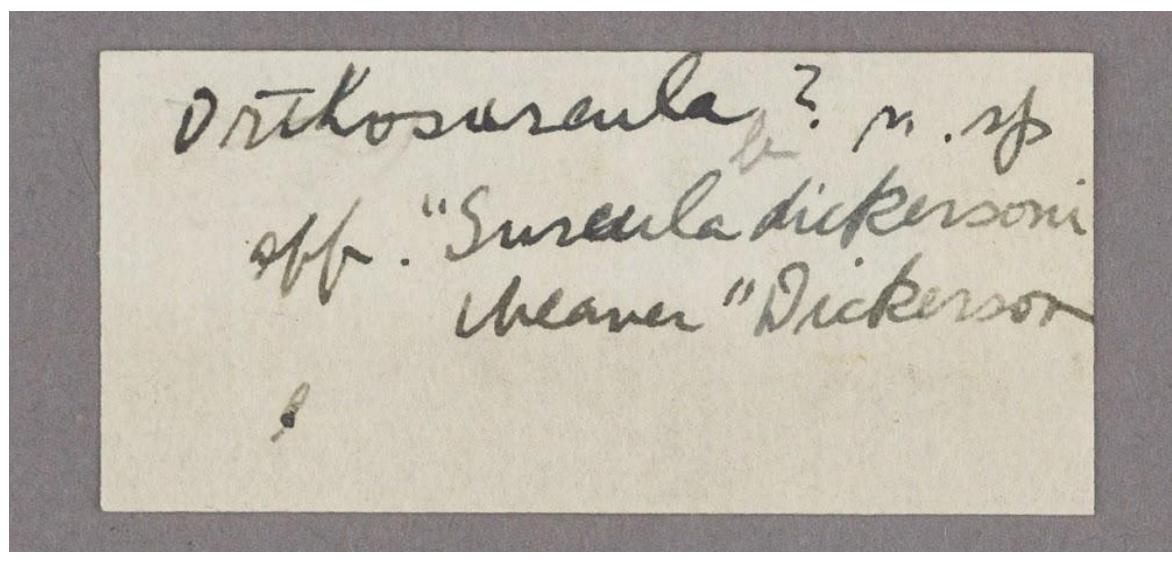
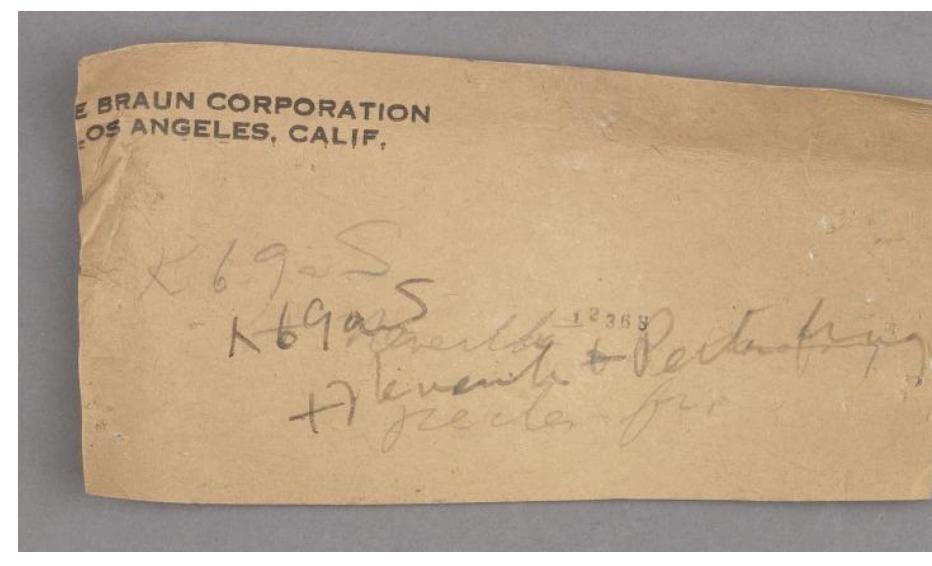


NMNH Paleo EPICC Stats

- 44,137 specimen lots imaged and databased
- 36,007 images of standard orientations
- 27,246 labels images and transcribed
- >220,000 specimens captured



EPICC is a partnership of nine natural history museums, united to digitize marine invertebrate fossils found in the eastern Pacific. The bivalves, gastropods, crabs, echinoids and other invertebrates found in our collections span the last 66 million years of Earth's history. During the course of the grant we will make 1.6 million specimen records available online through digital data and photographs. We are funded through the National Science Foundation's Advancing Digitization of Biological Collections program and affiliated with Integrated Digitized Biocollections (iDigBio). <https://epicc.berkeley.edu/>



1. Image Capture

With transcription, original data is processed into a standardized record. Imaging all labels creates a record of that original data. Metadata recorded in each specimen and label image was used to create skeleton records in the collections information system (EMu) and to initiate label transcription.

filename = catalog number captured from barcode
iptc:instructions = drawer location captured from drawer barcode
iptc:jobid = LocalityID in standardized format, found on label or specimen
(This process required creative use of iptc fields because the data did not map cleanly and is used for workflow only)

Prefix	Barcode	Event/ Site	Object ID	Drawer Location	Drawer Number	Drawer Location ID
PAL	USGS LOC 190			NHB - East Wing - Basement	EB1B - Row 08 - Column 05 - Lower Q	
PAL	USGS LOC 412			NHB - East Wing - Basement	EB1B - Row 08 - Column 02 - Upper Q	
PAL	USGS LOC 3548			NHB - East Wing - Basement	EB1B - Row 08 - Column 09 - Lower Q	
PAL	USGS LOC 1749			NHB - East Wing - Basement	EB1B - Row 08 - Column 04 - Upper Q	
PAL	USGS LOC 4287			NHB - East Wing - Basement	EB1B - Row 08 - Column 02 - Lower Q	
PAL	USGS LOC 4105			NHB - East Wing - Basement	EB1B - Row 08 - Column 02 - Lower Q	
PAL	USGS LOC 4089			NHB - East Wing - Basement	EB1B - Row 08 - Column 03 - Lower Q	
PAL	USGS LOC 4102			NHB - East Wing - Basement	EB1B - Row 08 - Column 03 - Lower Q	

3. Merging and Parsing with Python

We developed a Python script to merge the skeleton records and transcription records and to parse any identifiable data from "Additional Information" into more specific fields.

A benefit to creating the skeleton record from the image capture and then merging with transcription is that even specimens without labels will gain data.

2. Label Transcription

Transcription was completed in weekly batches creating a total of 21 datasets. It was important to create a workflow that would establish clear transcription guidelines to ensure consistency across the datasets and to account for the variety in data quality, quantity, and clarity as seen in the above label images.

Transcription Fields

1	Filename	Pulled from image file
2	Taxonomy	verbatim taxonomy string
3	Collector	
4	Geography	town, county, state
5	Location Description	
6	Additional Information	Any remaining information, labeled if possible (e.g. "Date: Aug 30-1917")

4. OpenRefine: Standardizing and Enhancing the Data

Each dataset is processed in OpenRefine, a powerful, open source data wrangling tool. With this tool we can complete data standardization and enhancement efficiently across these large datasets.

Standardization: Data from transcribed labels can vary in spelling, organization, and information. Using tools within OpenRefine such as clustering and facetting, patterns emerge enabling the cleaning and standardizing of the spelling of names, localities, and taxonomic names.

Using General Refine Expression Language (GREL) and Jython, mass standardization is made easier and steps can be repeated across datasets.

row value	value.replace('M.N.B.', 'M.N. Bramlett')
27 - W.P.	W.P. Woodring
116 - Kew	Kew
145 - W.M.B.	W.M.B. Woodring, M.N. Bramlett
176 - W.P.W.	W.P.W. Woodring, M.N. Bramlett
212 - M.B.	M.B. Woodring, M.N. Bramlett
265 - W.P.	W.P. Woodring, M.N. Bramlett

Using the GREL expression, only certain names were changed in the dataset.

OpenRefine will filter all of the data from a single column and analyse patterns in the data to suggest changes.

Enhancement: Data transcribed from labels can be lacking in context or information. OpenRefine provides tools for segmenting related data (i.e., collector name, taxonomic name, locality number) or finding recurring patterns. We can then use this information to reference outside resources such as ledgers, field books, publications, and records from other institutions for further data.

6 matching rows (1381 total)		
Show as: rows	records	Show: 5 10 25 50 rows
All	1299 ✓ LocNumber	Collector Geography Location Description Geologic Age Additional
All	1300 ✓ LocNumber	Collector Geography Location Description Geologic Age Additional
All	1301 ✓ LocNumber	Collector Geography Location Description Geologic Age Additional
All	1302 ✓ LocNumber	Collector Geography Location Description Geologic Age Additional
All	1303 ✓ LocNumber	Collector Geography Location Description Geologic Age Additional
All	1304 ✓ LocNumber	Collector Geography Location Description Geologic Age Additional

This process enables more information about the locality and stratigraphy can be included with the transcribed label data to form a more complete occurrence record. It also provides context for information that would otherwise lack meaning.

An essential step in this workflow includes documentation of enhancements and interpretations within the record. We use notes and verbatim fields to maintain transparency and long-term understanding of the data.

5. Backfilling with Python

As data was processed we maintained separate master files for all persons and geologic age data in order to track variation and expand the data based on existing records.

After OpenRefine, each dataset is processed with a second Python script that fills in higher level geologic age information, fields documenting updates to this information, and links to existing records for all persons.

```
FP_GEOAGE = 'geologic_age2.csv'
FP_PARTIES = 'parties.csv'
FP_INPUT = 'NMNH-PALEO-20180126_AC.openrefine.csv'
FP_OUTPUT = os.path.basename(FP_INPUT).rsplit('_', 1)[0] + '_emuprep.csv'

# def standardize(val):
#     """Standardizes the format of the string to improve matching"""
#     return val.strip().replace(' ', '').lower()

# def match_ages(records, errors=None):
#     """Matches new data to geologic ages already in EMU"""
#     if errors is None:
#         errors = []
#     # Read geologic age info from geologic_age2.csv
#     ages = {}
#     with open(FP_GEOAGE, 'r', encoding='utf-8') as gainput:
#         rows = csv.reader(gainput, dialect='excel')
#         keys = next(rows)
#         for row in rows:
#             ages[row[0]] = row[1]
#     for record in records:
#         if record['Geologic Age'] == 'Unknown':
#             record['Geologic Age'] = ages.get(record['Locality'], 'Unknown')
#     return records, errors
```

6. Ingest into EMu

Once the data is fully processed it is ingested into EMu, expanding the original skeleton records into more robust and standardized occurrence records with thorough documentation.

Pre	Catalog	Event/ Site	Object	Details	Collector	Vers	Comments	(Identification Details)	Note	(Note Details)	System	Geologic	Series	(Geo)
PAL	637989	Woodring	: USGS CENO LOC 12956	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Nevertia reclusiana (Deshayes, 1850)				EMU record was cr.	Neogene				
PAL	637990	Woodring	: USGS CENO LOC 12957	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Zirfaea cf. gibbi Tryon (1888)				EMU record was cr.	Neogene				
PAL	637991	Woodring	: USGS CENO LOC 12957	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Zirfaea cf. gibbi Tryon (1888)				EMU record was cr.	Neogene				
PAL	637992	Woodring	: USGS CENO LOC 12957	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Zirfaea cf. gibbi Tryon (1888)				EMU record was cr.	Neogene				
PAL	637993	Woodring	: USGS CENO LOC 12958	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Zirfaea cf. gibbi Tryon (1888)				EMU record was cr.	Neogene				
PAL	637994	Woodring	: USGS CENO LOC 12958	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Zirfaea cf. gibbi Tryon (1888)				EMU record was cr.	Neogene				
PAL	637995	Woodring	: USGS CENO LOC 12958	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Zirfaea cf. gibbi Tryon (1888)				EMU record was cr.	Neogene				
PAL	637996	Woodring	: USGS CENO LOC 12958	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Zirfaea cf. gibbi Tryon (1888)				EMU record was cr.	Neogene				
PAL	637997	Woodring	: USGS CENO LOC 12958	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Zirfaea cf. gibbi Tryon (1888)				EMU record was cr.	Neogene				
PAL	637998	Woodring	: USGS CENO LOC 12958	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Zirfaea cf. gibbi Tryon (1888)				EMU record was cr.	Neogene				
PAL	637999	Woodring	: USGS CENO LOC 12958	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Zirfaea cf. gibbi Tryon (1888)				EMU record was cr.	Neogene				
PAL	638000	Woodring	: USGS CENO LOC 12958	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Pseudocardium cf. macrostoma (Deshayes, 1850)				EMU record was cr.	Neogene				
PAL	638001	R. Stewart	: USGS CENO LOC 12958	USA: California: Kings County : Kettleman Hills: Aerial Gap Quadrangle Kettleman Hills, Aerial Gap Quadrangle N end of South Woodring	ID transcribed from label and not verified. Verbatim taxonomy: Pseudocardium cf. macrostoma (Deshayes, 1850)				EMU record was cr.					