

The Materials Experiment Knowledge Graph: Supporting Information

Michael J. Statt^{1,*}, Brian A. Rohr^{1,*}, Dan Guevarra^{2,3}, Ja’Nya Breedon³, Santosh K. Suram⁴, and John M. Gregoire^{2,3,*}

¹Modelyst LLC, Palo Alto, CA 94306

²Division of Engineering and Applied Science, California Institute of Technology, Pasadena, CA 91125

³Liquid Sunlight Alliance, California Institute of Technology, Pasadena, CA 91125

⁴Toyota Research Institute, Los Altos, CA 94022

*E-mail: brian.rohr@modelyst.io, michael.statt@modelyst.io, gregoire@caltech.edu

Node and edge types

| Entity Type | Count |
|-----------------|------------|
| Analysis | 3,980,832 |
| AnalysisDetails | 7 |
| Collection | 3,612 |
| Element | 46 |
| Process | 23,591 |
| ProcessData | 5,425,961 |
| ProcessDetail | 13,434 |
| Sample | 12,159,999 |
| SampleProcess | 30,656,471 |
| pH | 15 |

Table S1. Summary of entity count in the Neo4j graph database. Each row represents the count of a unique entity type.

| Relationship Type | Count |
|--------------------------|------------|
| ANALYSIS_DETAILS | 3,980,832 |
| CONTAINS | 3,527,845 |
| CollectionSample | 12,159,990 |
| NEXT | 18,590,083 |
| PH | 6,267 |
| PROCESS | 30,656,471 |
| PROCESS_DETAIL | 23,488 |
| ProcessDataAnalysis | 4,050,703 |
| SAMPLE | 3,065,6471 |
| SampleProcessProcessData | 7,777,908 |

Table S2. Summary of relationship count in the Neo4j graph database. Each row represents the count of a unique relationship type.

Computational Methods

The computational resources used to execute the use cases are as follows:

The extract, transform, load (ETL) process was carried out using a python library called DBgen (<https://github.com/modelyst/dbgen>), which was specifically designed to instantiate complicated, scientific data pipelines. PostgreSQL (<https://www.postgresql.org/>) was used to create the SQL database, and the Neo4j community edition (<https://neo4j.com/>) was used to create the graph database. The process of migrating the data from the SQL database to the graph database was done using a python library

Table S3. Resources used for each task

| Task | Resource Used |
|---|---|
| Extracting, transforming, and loading (ETL) the filesystem data into the SQL database | DBgen |
| Hosting the SQL database | PostgreSQL, AWS EC2, Docker |
| Migrating the SQL database to the graph database | PG4J |
| Hosting the graph database | Neo4j, AWS EC2, Docker |
| Querying graph database and processing data for design of experiments use cases | Python, Jupyter notebook server (local) |

17 called PG4J (<https://github.com/modelyst/pg4j>), which is capable of migrating any PostgreSQL database to Neo4j. For query
18 timing, the SQL database and the graph database were run in docker containers on AWS EC2. Specifically, the EC2 instance
19 was a t2.xlarge, and the docker images were postgres:14 and neo4j:5.5 for the SQL and graph databases, respectively. Cypher
20 queries and data processing for the design of experiments use cases were executed in Jupyter notebooks running on a local
21 JupyterHub server (Intel i9-11900K, 64 GB RAM). The computational methods are summarized in table S3.

22 Design of experiments use case

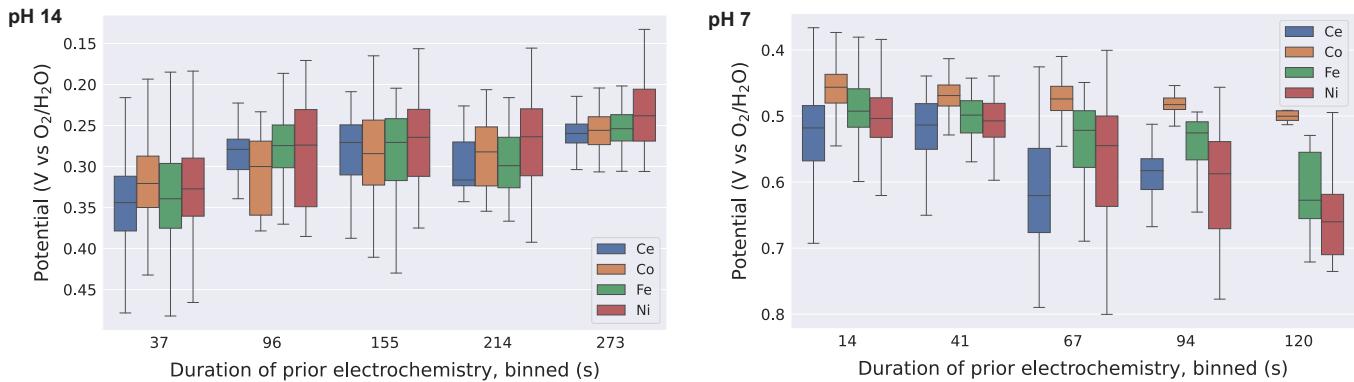


Figure S1. A summary of OER activity with 4982 measurements in pH 14 electrolyte (left) and 6524 measurements in pH 7 electrolyte (right). The catalysts with at least 70% concentration of the given element (Ce, Co, Fe, or Ni) are then grouped by 5 bins of total duration of electrochemical operation prior to the activity measurement. The catalyst overpotential at 3 mA/cm^2 is shown as an inverted vertical axis so that higher activity is shown as higher position in the figure. The consistent upward trend with increasing duration demonstrates a universal condition in pH 13 electrolyte, which is not observed in pH 7 electrolyte.

Table S4. Candidate composition spaces proposed by automated design of experiment for measurement in pH 3 electrolyte, sorted by 5th percentile overpotential at 10 mA/cm² from pH 7 measurement.

| Composition space | Eta (V) 3 mA/cm ² | Eta (V) 10 mA/cm ² | plate ID | sample count |
|-------------------|------------------------------|-------------------------------|----------|--------------|
| Co-Fe-La-Ni | 0.118 | 0.142 | 1740 | 1365.0 |
| Co-Fe-La-Ni | 0.118 | 0.142 | 1723 | 1365.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 1749 | 30.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 1751 | 30.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 1754 | 30.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 1755 | 30.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 1756 | 30.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 1762 | 388.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 1763 | 388.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 1774 | 30.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 2486 | 388.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 2487 | 388.0 |
| Ce-Co-Fe-Ni | 0.191 | 0.175 | 2488 | 388.0 |
| Ce-Co-Fe-La-Ni | 0.163 | 0.184 | 1757 | 1683.0 |
| Ce-Co-La-Ni | 0.435 | 0.483 | 1721 | 1398.0 |
| Ce-Co-La-Ni | 0.435 | 0.483 | 1720 | 1398.0 |
| Ce-Co-La-Ni | 0.435 | 0.483 | 2369 | 30.0 |
| Ce-Co-La-Ni | 0.435 | 0.483 | 2367 | 30.0 |
| Ce-Co-La-Ni | 0.435 | 0.483 | 1722 | 1398.0 |
| Ce-Co-La-Ni | 0.435 | 0.483 | 1750 | 30.0 |
| Ce-Co-La-Ni | 0.435 | 0.483 | 1719 | 1398.0 |
| Ce-Co-La-Ni | 0.435 | 0.483 | 1829 | 1398.0 |
| Co-Cu-Fe-Mn-Sn-Ta | 0.556 | 0.631 | 3673 | 3.0 |
| Co-Cu-Fe-Mn-Sn-Ta | 0.556 | 0.631 | 3859 | 63.0 |
| Co-Cu-Fe-Mn-Sn-Ta | 0.556 | 0.631 | 3865 | 63.0 |
| Co-Cu-Fe-Mn-Sn-Ta | 0.556 | 0.631 | 3866 | 63.0 |
| Co-Cu-Fe-Mn-Sn-Ta | 0.556 | 0.631 | 3867 | 63.0 |
| Co-Cu-Fe-Mn-Sn-Ta | 0.556 | 0.631 | 3870 | 63.0 |
| Ce-Co-Ni-Zn | 0.521 | 0.634 | 2532 | 1597.0 |

23 **Creation of Database Subsets**

24 To investigate the scaling of query times in both the SQL and graph databases, we created three smaller versions of the original
25 database. The first database fragment was created by removing the last half of the rows in the sample-process table, ordered by
26 their process timestamps. We then deleted all rows in other tables that were no longer linked to a sample-process. This process
27 was repeated two more times to create two additional database fragments, with 3/4 and 7/8 of the rows in the sample-process
28 table deleted. Each fragment was migrated to Neo4j using the tools described above, resulting in a series of MPS-style and
29 MEKG-style databases that share the same information and contain 1/8, 1/4, and 1/2 of the number of Sample-Processes in the
30 full MPS and MEKG databases.

31 **Query 4 in Cypher and SQL**

32 Below is the code for Query 4 in Cypher:

```
33   MATCH
34    path=(sp1:SampleProcess)-[:NEXT]->(sp2)-[:NEXT]->(sp3)-[:NEXT]->(sp4)-[:NEXT]->(sp5),
35    (a1:Analysis)<--(pda1:ProcessData)<--(sp1)-->(p1:Process)-->(pd1:ProcessDetail),
36    (a2:Analysis)<--(pda2:ProcessData)<--(sp2)-->(p2:Process)-->(pd2:ProcessDetail),
37    (a3:Analysis)<--(pda3:ProcessData)<--(sp3)-->(p3:Process)-->(pd3:ProcessDetail),
38    (a4:Analysis)<--(pda4:ProcessData)<--(sp4)-->(p4:Process)-->(pd4:ProcessDetail),
39    (a5:Analysis)<--(pda5:ProcessData)<--(sp5)-->(p5:Process)-->(pd5:ProcessDetail)
40 WHERE
41 pd1.technique STARTS WITH 'CA'
42 AND pd2.technique STARTS WITH 'CA'
43 AND pd3.technique STARTS WITH 'CA'
44 AND pd4.technique STARTS WITH 'CA'
45 AND pd5.technique STARTS WITH 'CV'
46 AND a5.name = 'CV_FOMS_standard'
47 AND apoc.convert.fromJsonMap(a1.output) ['I.A_ave'] > 1e-7
48 AND apoc.convert.fromJsonMap(a2.output) ['I.A_ave'] > 1e-8
49 AND apoc.convert.fromJsonMap(a3.output) ['I.A_ave'] > 1e-9
50 AND apoc.convert.fromJsonMap(a4.output) ['I.A_ave'] > 1e-10
51 AND apoc.convert.fromJsonMap(a5.output) ['I.A_max'] > 1e-6
52 AND apoc.convert.fromJsonMap(pd1.parameters) ['electrolyte'] CONTAINS 'NaOH'
53 AND apoc.convert.fromJsonMap(pd2.parameters) ['electrolyte'] CONTAINS 'NaOH'
54 AND apoc.convert.fromJsonMap(pd3.parameters) ['electrolyte'] CONTAINS 'NaOH'
55 AND apoc.convert.fromJsonMap(pd4.parameters) ['electrolyte'] CONTAINS 'NaOH'
56 AND apoc.convert.fromJsonMap(pd5.parameters) ['electrolyte'] CONTAINS 'NaOH'
57 RETURN count(path)
```

58 Below is the code for Query 4 in SQL:

```
59 with your_table as (
60 select
61   sp.sample_id,
62   p1."timestamp",
63   p1."ordering",
64   pd1.technique,
65   pd1.parameters,
66   a."name",
67   a."output"
68 from
69   sample_process sp
70 join process p1 on
71   sp.process_id = p1.id
72 left join process_detail pd1 on
73   p1.process_detail_id = pd1.id
74 left join sample_process_process_data sppd on
75   sppd.sample_process_id = sp.id
76 left join process_data pd on
77   sppd.process_data_id = pd.id
78 left join process_data_analysis pda on
79   pda.process_data_id = pd.id
80 left join analysis a on
81   pda.analysis_id = a.id
82 where
83   sp.sample_id in (
```

```

84  select
85    sp3.sample_id
86  from
87    sample_process sp3
88  join process p3 on
89    sp3.process_id = p3.id
90  join process_detail pd3 on
91    p3.process_detail_id = pd3.id
92  where
93    pd3.technique like 'CV%'
94  )
95  and sp.sample_id in (
96  select
97    sp2.sample_id
98  from
99    sample_process sp2
100 join process p2 on
101   sp2.process_id = p2.id
102 join process_detail pd2 on
103   p2.process_detail_id = pd2.id
104 where
105   pd2.technique like 'CA%'
106 group by
107   sp2.sample_id
108 having
109   count(*) >= 4
110 ),
111 filtered_labels as (
112 select
113   sample_id
114 from
115  your_table
116 group by
117   sample_id
118 having
119   COUNT(*) >= 5
120 ),
121 sequenced_data as (
122 select
123   t1.sample_id,
124   t1.Timestamp,
125   t1.ordering,
126   t1.technique,
127   t1.parameters,
128   t1."name",
129   t1."output",
130   row_number() over (partition by t1.sample_id
131 order by
132   t1.Timestamp,
133   t1.ordering) as RowNum
134 from
135  your_table t1
136 inner join
137  filtered_labels fl on
138   t1.sample_id = fl.sample_id

```

```

139 ),
140 json_agg_data as (
141 select
142   sd1.sample_id,
143   json_agg(sd1.technique) over (partition by sd1.sample_id
144 order by
145   sd1.RowNum rows between current row and 4 following) as technique_seq,
146   json_agg(sd1.parameters) over (partition by sd1.sample_id
147 order by
148   sd1.RowNum rows between current row and 4 following) as parameters_seq,
149   json_agg(sd1.name) over (partition by sd1.sample_id
150 order by
151   sd1.RowNum rows between current row and 4 following) as name_seq,
152   json_agg(sd1.output) over (partition by sd1.sample_id
153 order by
154   sd1.RowNum rows between current row and 4 following) as output_seq,
155   COUNT(*) over (partition by sd1.sample_id
156 order by
157   sd1.RowNum rows between current row and 4 following) as technique_seq_count
158 from
159 sequenced_data sd1
160 )
161 select
162   count(*)
163   -- sample_id
164   -- technique_seq,
165   -- parameters_seq,
166   -- name_seq,
167   -- output_seq
168 from
169 json_agg_data
170 where
171   technique_seq_count = 5
172   and technique_seq->>0 like 'CA%'
173   and technique_seq->>1 like 'CA%'
174   and technique_seq->>2 like 'CA%'
175   and technique_seq->>3 like 'CA%'
176   and technique_seq->>4 like 'CV%'
177   and name_seq->>4 = 'CV_FOMS_standard'
178   and (output_seq->0->>'I.A_ave')::float > 1e-7
179   and (output_seq->1->>'I.A_ave')::float > 1e-8
180   and (output_seq->2->>'I.A_ave')::float > 1e-9
181   and (output_seq->3->>'I.A_ave')::float > 1e-10
182   and (output_seq->4->>'I.A_max')::float > 1e-6
183   and parameters_seq->0->>'electrolyte' like '%NaOH%'
184   and parameters_seq->1->>'electrolyte' like '%NaOH%'
185   and parameters_seq->2->>'electrolyte' like '%NaOH%'
186   and parameters_seq->3->>'electrolyte' like '%NaOH%'
187   and parameters_seq->4->>'electrolyte' like '%NaOH%'
188

```