1    # Ontology of RNA Sequencing (ORNASEQ) and its application

2

3    Stephen A Fisher

4    Biology Department

5    School of Arts and Sciences

6    University of Pennsylvania

7    433 S. University Avenue

8    Philadelphia, PA 19104

9    safisher@upenn.edu

10

11    Junhyong Kim*

12    Biology Department

13    School of Arts and Sciences

14    University of Pennsylvania

15    433 S. University Avenue

16    Philadelphia, PA 19104

17    junhyong@upenn.edu

18

19    * Corresponding author

20

## Abstract

**Background:** Next-generation RNA sequencing is a rapidly developing technology with complex procedures encompassing different experimental modalities. As the technology evolves and its use expand, so does the need to capture the data provenance from these sequencing studies and the need to create new tools to manage and manipulate these provenance stores.

**Results:** Here we used the Ontology for Biomedical Investigations (OBI) and many other ontologies from the Open Biological and Biomedical Ontology (OBO) Foundry as a framework from which to create an application ontology (ORNASEQ: Ontology of RNA sequencing) to capture data provenance for next-generation RNA sequencing studies. Additionally, we provide an extensive real-life sample provenance data set for use in developing new provenance tools and additional sequencing data models.

**Conclusions:** The Ontology of RNA Sequencing (ORNASEQ) provides core terms for use in building data models to capture the provenance from next-generation RNA sequencing studies. The supplied sample provenance data also exemplifies many of the complexities of RNA sequencing studies and underscores the need for potent workflow management systems.

## Keywords

Ontology, RNAseq, PROV-XML

Fisher, Kim    2

43

## Background

45   Until recently the cost of performing next-generation RNA sequencing (RNAseq)

46   experiments limited the amount of data generated by a single lab and managing and

47   properly documenting a few experiments was not fundamentally challenging.

48   However, as sequencing costs have dropped, research groups are now running

49   hundreds, thousands or even tens of thousands of RNAseq experiments, creating a

50   need to systematically document experimental and informatics details and track

51   provenance of the final published or publicly released datasets. RNAseq has also

52   begun making its way into medical diagnostics, where data provenance is a

53   necessity for quality assurance and regulatory compliance. Tracking the data

54   provenance for hundreds or thousands of sequencing experiments in either a

55   research or medical setting requires data models and structures that must be put

56   into place to capture the necessary information at all stages of a sequencing

57   experiment and it's not always obvious what information is necessary. While there

58   are numerous platforms and pipelines to analyze RNAseq data, there are limited

59   data models or ontologies that could be applied to successfully capture the details of

60   an RNAseq experiment [1–6].

61

62   Within a single next-generation sequencing experiment there is a dizzying amount

63   of information that must be captured throughout the often-complex experimental

64   procedures and post-sequencing informatics analyses. The problem is further

65    complicated by the number of researchers or technicians who might be involved in a

66    single sequencing experiment whose roles are interspersed in irregular patterns;

67    complexity of biological specimens, their origin and experimental designs; and, the

68    frequent disconnect between the biologists running the experiments and the

69    bioinformaticists analyzing the data. Tracking data provenance spanning

70    experiment procedures recorded in lab notebooks belonging to multiple biologists

71    and computer log files residing in a series of cryptic directories on a file system

72    quickly becomes an intractable problem. These challenges suggest a need for a

73    comprehensive next generation sequencing provenance system. Data provenance

74    requires data models, provenance models, and supporting infrastructure. Here, we

75    focus on the first part of data models for RNA sequencing experiments and describe

76    an Ontology for RNA Sequences (ORNASEQ).  In addition, we provide a large next-

77    generation sequencing use case from an active RNAseq workflow, using the PROV-

78    XML database format for the community to use as an example dataset for

79    development of provenance models and tools.

## Ontology for RNA sequencing

81    The Ontology for RNA Sequencing (ORNASEQ) is an application ontology based

82    largely on the Ontology for Biomedical Investigations (OBI)[1], using the principles

83    of OBO Foundry[7]. Specifically, ORNASEQ contains 162 terms, 117 of the terms are

84    from 16 existing ontologies, with 28 new terms having now been added to OBI and

85    17 terms being defined directly in ORNASEQ (see Table 1). ORNASEQ is designed to

86    annotate RNA-based next-generation sequencing, although much of ORNASEQ

Fisher, Kim    4

87   would also apply to DNA-based next-generation sequencing. The ontology was

88   designed, in part, through efforts to track the data provenance of thousands of

89   RNAseq samples collected by the NIH Common Fund Single Cell Analysis Program-

90   Transcriptome (SCAP-T) program[8]. Terms included in the ontology cover pre-

91   sequencing preparations and primary post-sequencing data analysis.

92

| Ontology | Number Terms |
|---|---|
| **BFO**: Basic Formal Ontology<br><br>*http://purl.obolibrary.org/obo/bfo.owl* | 10 |
| **CHEBI**: Chemical Entities of Biological Interest<br><br>*http://purl.obolibrary.org/obo/chebi.owl* | 2 |
| **CL**: Cell Ontology<br><br>*http://purl.obolibrary.org/obo/cl.owl* | 6 |
| **EFO**: Experimental Factor Ontology<br><br>*http://www.ebi.ac.uk/efo/efo.owl* | 6 |
| **GENEPIO**: Genomic Epidemiology Ontology<br><br>*http://purl.obolibrary.org/obo/genepio.owl* | 4 |
| **GO**: Gene Ontology<br><br>*http://purl.obolibrary.org/obo/go.owl* | 5 |
| **IAO**: Information Artifact Ontology<br><br>*http://purl.obolibrary.org/obo/iao.owl* | 14 |
| **NCBITaxon**: NCBI Organismal Classification<br><br>*http://purl.obolibrary.org/obo/ncbitaxon.owl* | 1 |
| **NCIT**: NCI Thesaurus OBO Edition<br><br>*http://purl.obolibrary.org/obo/ncit.owl* | 6 |
| **OBI**: Ontology for Biomedical Investigations<br><br>*http://purl.obolibrary.org/obo/obi.owl* | 79 |

| | |
|---|---|
| **OBIws**: OBI web service, development version<br><br>*http://purl.obolibrary.org/obo/obi/webService.owl* | 5 |
| **OGMS**: Ontology for General Medical Science<br><br>*http://purl.obolibrary.org/obo/ogms.owl* | 1 |
| **OMIABIS**: Ontologized MIABIS<br><br>*http://purl.obolibrary.org/obo/omiabis.owl* | 2 |
| **SO**: Sequence Types and Features<br><br>*http://purl.obolibrary.org/obo/so.owl* | 1 |
| **TAXRANK**: Taxonomic rank vocabulary<br><br>*http://purl.obolibrary.org/obo/taxrank.owl* | 2 |
| **UBERON**: Uberon multi-species anatomy ontology<br><br>*http://purl.obolibrary.org/obo/uberon.owl* | 1 |

93   Table 1: The set of external ontologies included in ORNASEQ.

94

95   Pre-Sequencing Provenance

96   Knowing what happens to a data sample prior to sequencing is essential to

97   understanding the analyzed data. The provenance surrounding the preparation of

98   sequencing data can prove invaluable to diagnosing aberrant results. These

99   protocols are often revised over time and as hardware and reagents evolve and a

100   multi-year study will likely include many versions of protocols with varying degrees

101   of differential change. Even in the controlled context of medical diagnostics, changes

102    are inevitable, for example, when a reagent becomes discontinued or hardware are

103    upgraded.

104

105    When capturing experimental lab data provenance, it is difficult to know what to

106    capture and in what format (e.g. as fields in a database, as a Word or PDF document,

107    or even as a reference to an entry in an electronic notebook). As sequencing

108    preparation protocols are often distributed as PDF or Word documents, trying to

109    track changes across multiple such documents quickly becomes a tedious process

110    that is difficult to automate and nearly impossible to query for specific questions

111    (e.g., what version of sequencing chemistry was used for what samples). Conversely,

112    because the protocols involve many incremental steps and details that are inter-

113    related in a complicated manner, it is difficult to convert each protocol into a fine-

114    grained structured knowledge model suitable for standard DBMS. Similar to the

115    scheme implemented by the Genomic Standards Consortium [3], we propose

116    experimental provenance be captured in a hybrid fashion including both detailed

117    protocol files and hardcoded fields tracking fundamental datum external to the

118    protocol files. In this RNAseq data provenance schema, complete protocol

119    definitions are stored in the provenance database as files (typically PDF or Word

120    documents). Critical or commonly changed features of the protocol are additionally

121    captured in the database schema. While not optimal, this approach preserves the

122    fine detail required to replicate an experimental procedure, while allowing for

123    structured query of the main features.

124

Fisher, Kim    8

## 125    Post-Sequencing Provenance

126    When RNAseq data comes off the sequencer it is typically converted into fastq files

127    by a proprietary, sequencer-specific program. A fastq file is a text file containing

128    nucleotide sequences [9]. In the case of RNAseq, fastq files contain nucleotide

129    sequences representing the RNA molecules from a biological sample. Any one RNA

130    molecule might be represented by tens to hundreds of thousands of sometimes

131    duplicate nucleotide sequences in the fastq file, with each nucleotide sequence

132    referred to as a "read" (i.e. a read out of the biological sequence). The fastq files are

133    run through what is called a "primary analysis pipeline", which may include any

134    number of steps such as generating quality control metrics, removing contaminant

135    sequences from one or both ends of the reads, removing low quality sub-sequences

136    from reads, removing duplicate reads, and aligning reads to a pre-defined reference

137    library. Primary analysis pipelines and more generally RNA sequencing experiments

138    often culminate in the generation of gene-, transcript-, or exonic-counts of the

139    number of reads associated with each of these categories. The various steps will

140    mostly remain consistent within a research group, a project or a medical diagnostic

141    test. However, programs evolve, algorithmic bugs are fixed, and reference libraries

142    are refined. As with wet lab procedures, it is common and often necessary for

143    primary analysis pipelines to change by varying degrees over time. Since the data

144    and processes for the informatics steps are already machine readable, capturing

145    provenance for RNAseq analysis pipeline is a more obvious task than in a wet lab

146    context. However, relevant provenance information is typically realized as a set of

147    log files spread across a series of programs, each with specific directory structures.

148  Using log files to track provenance data also quickly breaks down as programs

149  change or pipelines evolve and often requires complex programming to perform

150  simple queries across data samples. As with tracking wet lab data, provenance from

151  sample processing pipelines needs to be rigorously captured and systematically

152  stored in a central database. Addressing these challenges require incorporation of a

153  well-defined and use-case oriented ontology for the provenance objects, which we

154  provide with ORNASEQ.

## 155  Sample Data

156  ORNASEQ is meant to provide core terms used to track the provenance of RNAseq

157  datasets. However, any particular experiment or RNAseq use case will require a

158  multitude of additional terms. Here we provide a dataset containing curated and

159  modified data provenance from 1,347 next-generation sequencing samples. The

160  dataset contains 93 data fields, using ORNASEQ terms as appropriate. Each sample

161  was processed with one of four different versions of a primary analysis pipeline and

162  there was from one to three variants (sub-versions) of each pipeline version.

163  Specifically, samples were processed either as "single-end" or "paired-end" and

164  aligned with either STAR[10] or Bowtie[11, 12], ultimately leading to nine possible

165  pipeline variants. The dataset is provided as PROV-XML (see Additional file 2), with

166  the data summarized in an Excel table (see Additional file 1). Real world use cases

167  usually include messy data. Occasionally pipelines are run incorrectly either

168  through intentional operator actions or error. It's also quite common to have

169  missing or incomplete data provenance again through user actions (e.g. data was

170   pulled from a source that lacked sufficient provenance), programming errors,

171   computer issues, etc. The dataset provided here intentionally includes both

172   incorrect and missing data.

173

174   Primary Analysis Pipeline Stages

175   The data provided describes an analysis pipeline with seven possible stages but

176   with each analysis only including five of the seven stages. Table 2 includes a subset

177   of data tracked for each pipeline stage and which stages might be used in a

178   particular pipeline version. For example, the stage HTSeq was only used in pipeline

179   version 1.0 while VERSE was only used in later pipeline versions. The stages are

180   very briefly described here and more complete descriptions can be found in the

181   PennSCAP-T Pipeline[13].

182         o   BLAST - a subset of samples was aligned to the Blast NR database,

183               using BLAST as a quality control check.

184         o   FastQC – FastQC is used as an additional quality control check.

185         o   TRIM – contaminant sequences were removed from reads.

186         o   BOWTIE – Bowtie was used to align reads to a reference genome.

187         o   STAR – STAR was used to align reads to a reference genome.

188         o   HTseq – HTSeq was used to assign aligned reads to genes.

189         o   VERSE – VERSE was used to assign aligned reads to genes.

190

191   In the accompanying dataset, samples were aligned with either Bowtie or STAR but

192   not both. Similarly, alignments were processed by either HTSeq or VERSE, but not

Fisher, Kim  11

193   both. The Excel table highlights which provenance terms are consistent within a

194   pipeline version. For example "BLAST.blastn.version" has pipeline version-specific

195   values. Provenance terms that might contain incorrect values are also denoted. For

196   example, the value for "STAR.star.version" in Pipeline version 2.0.1 might be

197   missing.

198

| STAGE | PARAMETER | VALUE | | | |
|---|---|---|---|---|---|
| **BLAST** | blastn version | 2.2.25 | 2.2.30 | 2.2.30 | 2.2.30 |
| | parseBlast.py version | 1.1 | 1.5 | 1.5 | 1.8 |
| **FASTQC** | fastqc version | 0.10.1 | 0.11.2 | 0.11.2 | 0.11.2 |
| **TRIM** | trimReads.py version | 0.6 | 0.6 | 0.6 | 0.7.2 |
| | removeN | 1 | 1 | 1 | 1 |
| | minLen | 30 | 20 | 20 | 20 |
| | phredThresh | - | 53 | 53 | 53 |
| | numAT | 30 | 26 | 26 | 26 |
| **BOWTIE** | bowtie version | 0.12.7 | 1.1.1 | | 2.2.9 |
| | minins | 250 | 150 | | 150 |
| | maxins | 450 | 600 | | 600 |
| **STAR** | star version | 2.3.0.1 | 2.4.0h1 | 2.4.0j | 2.5.1b |
| | samtools version | 1.1 | 1.1 | 1.2 | 1.2 |
| **HTSEQ** | htseq version | 0.5.4p5 | | | |
| **VERSE** | verse version | | v0.1.4 | v0.1.4 | v0.1.5 |
| | **Pipeline Version** | **1.0** | **2.0** | **2.0.1** | **2.1** |
| | **Number Samples** | **208** | **413** | **216** | **510** |

199    Table 2: Subset of pipeline-specific parameters included in sample dataset, across

200    pipeline stages and pipeline versions.

## Conclusion

202    Next generation sequencing (NGS) is a complex technology with multiple varied

203    steps. Capturing the provenance of NGS data requires complex systems and multi-

204    user collaborations. The schemas and ontology we define here offer a basic

205    framework that can be tuned and expended by researchers to the particulars of

206    individual studies, while providing basic commonalities across studies.

207

208    The dataset provided illustrates some of the complexities of the data provenance

209    from downstream processing of NGS data. These complexities will grow as the field

210    evolves. For example, there are now hundreds of variants of basic sequencing

211    protocol that are specific to particular biology applications (e.g., ATAC-seq; [14];

212    Drop-Seq;[15, 16]). Each of these involve variations in experimental and informatics

213    processes. It will become necessary to build workflow management systems and

214    smart clustering algorithms of the provenance[17] to help segment incomplete and

215    erroneous data. We propose that our large real-life dataset example will prove

216    useful in designing future new workflow systems and provenance models.


217    ## List of Abbreviations

218    **OBI:** Ontology for Biomedical Investigations

219    **ORNASEQ:** Ontology of RNA Sequencing

220    **OBO:** Open Biological and Biomedical Ontology

221    **RNAseq:** next-generation RNA sequencing

222    **SCAP-T:** NIH Common Fund Single Cell Analysis Program-Transcriptome

223    **NGS:** Next generation sequencing

## 224 Declarations

### 225 Ethics approval and consent to participate

226 not applicable

### 227 Consent for publication

228 not applicable

### 229 Availability of data and material

230 The ontology presented in the current study is available from GitHub,

231 http://doi.org/10.5281/zenodo.1311869. The sample dataset presented in the

232 current study is available as additional article files.

### 233 Competing interests

234 The authors declare that they have no competing interests.

### 235 Funding

236 This work was supported in part by NIH grants 5U01EB020954 and U01MH098953.

237 The NIH played no role in the results of this research effort or the text of this paper.

### 238 Authors' contributions

239 SF and JK contributed to the design of this study and preparation of the manuscript.

### 240 Acknowledgements

241 We are grateful to Christian Stoeckert and Daniel Berrios for valuable feedback and

242 assistance adding terms to the Ontology for Biomedical Investigations.

## References

243

244    1. Bandrowski A, Brinkman R, Brochhausen M, Brush MH, Bug B, Chibucos MC, et al.

245    The Ontology for Biomedical Investigations. PLoS One. 2016;11:e0154556.

246    doi:10.1371/journal.pone.0154556.

247    2. Society FGD. Minimum Information about a high-throughput SEQuencing

248    Experiment. http://fged.org/projects/minseqe/.

249    3. Yilmaz P, Kottmann R, Field D, Knight R, Cole JR, Amaral-Zettler L, et al. Minimum

250    information about a marker gene sequence (MIMARKS) and minimum information

251    about any (x) sequence (MIxS) specifications. Nat Biotechnol. 2011;29:415–20.

252    doi:10.1038/nbt.1823.

253    4. Sansone SA, Rocca-Serra P, Field D, Maguire E, Taylor C, Hofmann O, et al. Toward

254    interoperable bioscience data. Nat Genet. 2012;44:121–6. doi:10.1038/ng.1054.

255    5. Global Alliance for Genomics & Health. 2018. https://www.ga4gh.org.

256    6. Hong EL, Sloan CA, Chan ET, Davidson JM, Malladi VS, Strattan JS, et al. Principles

257    of metadata organization at the ENCODE data coordination center. Database

258    (Oxford). 2016;2016. doi:10.1093/database/baw001.

259    7. Smith B, Ashburner M, Rosse C, Bard J, Bug W, Ceusters W, et al. The OBO

260    Foundry: coordinated evolution of ontologies to support biomedical data

261    integration. Nat Biotechnol. 2007;25:1251–5. doi:10.1038/nbt1346.

262    8. Single Cell Analysis Program – Transcriptome Project. http://scap-t.org/.

263    9. Cock PJ, Fields CJ, Goto N, Heuer ML, Rice PM. The Sanger FASTQ file format for

264    sequences with quality scores, and the Solexa/Illumina FASTQ variants. Nucleic

265    Acids Res. 2010;38:1767–71. doi:10.1093/nar/gkp1137.

266    10. Dobin A, Davis C a, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR:

267    ultrafast universal RNA-seq aligner. Bioinformatics. 2013;29:15–21.

268    doi:10.1093/bioinformatics/bts635.

269    11. Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient

270    alignment of short DNA sequences to the human genome. Genome Biol.

271    2009;10:R25. doi:10.1186/gb-2009-10-3-r25.

272    12. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat

273    Methods. 2012;9:357–9. doi:10.1038/nmeth.1923.

274    13. Fisher  J. SA. K. PennSCAP-T Pipeline. 2018. https://github.com/safisher/ngs.

275    14. Buenrostro JD, Wu B, Chang HY, Greenleaf WJ. ATAC-seq: A Method for Assaying

276    Chromatin Accessibility Genome-Wide. Curr Protoc Mol Biol. 2015;109:21 29 1-9.

277    doi:10.1002/0471142727.mb2129s109.

278    15. Klein AM, Mazutis L, Akartuna I, Tallapragada N, Veres A, Li V, et al. Droplet

279    barcoding for single-cell transcriptomics applied to embryonic stem cells. Cell.

280    2015;161:1187–201. doi:10.1016/j.cell.2015.04.044.

281    16. Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, et al. Highly

282    Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter

283    Droplets. Cell. 2015;161:1202–14. doi:10.1016/j.cell.2015.05.002.

284    17. Alawini L.; Davidson, S.; Fisher, S.; Kim, J. A. C. Discovering Similar Workflows via

285 Provenance Clustering: a Case Study. In: Belhajjame K, Gehani A, Alper P, editors.

286 Provenance and Annotation of Data and Processes. London: Springer International

287 Publishing; 2018.

288

## Additional Files

289

290 Additional file 1

291 **Format:** Excel Workbook from Excel for Mac version 16 (.xlsx)

292 **Title**: A Large Curated RNAseq Metadata Dataset

293 **Description**: This file contains a curated dataset consisting of metadata collected

294 from 1347 next-generation sequencing samples. The file also contains a set of terms

295 that can be used to separate the data samples between nine primary analysis

296 pipelines.

297 Additional file 2

298 **Format**: XML file (.xml)

299 **Title**: PROV-XML Instantiation of a RNAseq Metadata Dataset

300 **Description**: This is an XML file that contains a PROV-XML representation of the

301 data included in "Additional file 1.xslx"; that is, an example dataset containing

302 metadata from 1347 RNAseq data samples that were processed with one of nine

303 primary analysis pipelines.