

1 MetroNome — a visual data exploration platform for integrating 2 human genotypic and phenotypic data across diseases

3 Christian Stolte^{1*}, Kevin Shi¹, Nina Lapchyk^{1*}, Nathaniel Novod¹, Avinash Abhyankar¹, Lyle W.
4 Ostrow², Hemali Phatnani^{1,3}, Toby Bloom^{1*}

5

6 Affiliations

7 1 New York Genome Center, New York, NY 10013

8 2 Department of Neurology, Johns Hopkins University School of Medicine,
9 Baltimore, MD 21205

10 3 Columbia University Motor Neuron Center, New York, NY 10032

11

12 * affiliated at time of project completion

13 Abstract

14 MetroNome is a web-based visual data exploration platform which integrates de-identified
15 genomic, transcriptomic, and phenotypic data sets. Users can define and compare cohorts
16 constructed from multimodal data and share the data and analyses with outside tools. MetroNome's
17 interactive visualization and analysis tools allow researchers to quickly form and explore novel
18 hypotheses. The deidentified data is linked back to the source biosample inventories in multiple
19 biobanks, enabling researchers to further investigate new ideas using the most relevant samples.

20 Introduction

21 Biomedical research is producing a wealth of genomic data, some of it public [1]; though much is
22 restricted to various consortia or project team members. The restrictions are often necessary to
23 comply with consents and regulatory policies. Analyses of complex multimodal patient-derived
24 data—such as genome sequencing, clinical and pathological measures, environmental factors, and
25 imaging—enables questions to be explored that would otherwise not be possible. For example, how
26 can the same chromatin remodeling genes be associated with autism, schizophrenia, bipolar
27 disorder, congenital heart defects, and digestive tract issues? To achieve adequate statistical power
28 for genomic research discoveries, different types of data — from different studies and diseases —
29 must be integrated while assuring regulatory compliance with patient confidentiality and data use
30 policies [2, 3, 4, 6]. The associated tools should be openly available and usable by a wide audience
31 with different levels of expertise in genomics and biostatistics, while still ensuring responsible use of
32 the data. Existing applications make it possible to integrate GWAS results with other data to
33 prioritize variants by phenotype [7] or browse available individual-level genotype and sequence data

34 associated with phenotypic features [8]. However, these tools generally lack the capability to then
35 generate lists of samples and subjects matching genomic and phenotypic criteria of interest, enabling
36 access to the underlying de-identified data and samples across multiple biorepositories.

37 MetroNome comprises a web-based suite of interactive data visualization tools, enabling users to
38 combine their own data with other relevant public datasets and explore the results via linked graphs
39 and diagrams. The system leverages human visuospatial cognitive abilities to reveal patterns and find
40 connections. Researchers can segregate results along multiple dimensions, such as tissue source,
41 quality control measures, whether specific variants are present, or any combination of ranges and
42 categories in phenotypic and demographic dimensions. This enables easy comparisons between user-
43 defined cohorts via the visualization modules. Normalized analyses, such as gene expression z-scores,
44 are calculated in real time, based on the user's search parameters, to allow side-by-side visual
45 comparisons which highlight critical differences between groups, e.g., cases and controls.

46

47 Background

48 MetroNome derives from two thus far distinct paths in data exploration. Genomic visualization
49 tools such as cBioPortal [14] provide domain-specific visualizations such as variant diagrams and the
50 oncoprint visualization. An earlier and orthogonal stream of work was the development of
51 commercial “OnLine Analytic Processing” (OLAP) tools, which date to the 1970's but gained
52 traction in the 1990's with tools such as Cognos (now IBM) [www.ibm.com/Cognos/Analytics/], and
53 Business Objects (now SAP) [www.businessobjects.com/]. These tools enabled multi-dimensional
54 analysis of data, and dynamic filtering of commercial data along individual dimensions – although
55 they were primarily tabular, not graphic. The second generation of OLAP, led by Tableau [
56 www.tableau.com] introduced graphics and dashboards. Dashboards provide multiple visualizations
57 on the same screen and the ability to filter on any one frame and propagate that filter to all other
58 frames on the dashboard. This dynamic queries technique initially arose as an alternative to SQL for
59 querying databases [5]. MetroNome unites these two directions: scientific and statistical
60 visualizations, combined with dynamic filtering and filter propagation.

61 Exploration of synthetic cohorts via multi-modal interactive data visualization

62 MetroNome presents genomic and transcriptomic data in the context of phenotypic attributes,
63 relying on customized linked visualizations to enable exploration.

64 Creating synthetic cohorts

65 We allow the user to combine and display data from multiple sources, based on phenotypic or
66 genotypic traits and user access rights to those datasets. We provide access to publicly available
67 reference datasets, such as 1000 genomes [12] and TCGA tumor data [13] for use as comparators. To
68 perform an analysis, the user starts by creating a query that selects data from one or multiple

69 sources. MetroNome's query page presents a series of dynamically linked drop-down menus that can
70 be combined into rules for selecting subjects. Rules for subject and sample criteria can include filters
71 for any information available in the datasets that the user has selected for inclusion. Genomic rules
72 can include the presence of variants in a list of genes or in a genomic region. Variants can be filtered
73 by their predicted protein-coding impact as calculated by SnpEff [9], or their association with
74 disease as recorded in ClinVar [10]. The search produces lists of subjects, samples, and variants that
75 match the selected rules.

76 Multi-modal visualization and refinement of cohorts

77 The query results display multiple types of data simultaneously in linked panels, such as whole
78 genome or exome variants, RNA expression, copy number variations, structural variants, and
79 phenotypic data, to facilitate visual exploration of associations among different data types. The
80 visualization controls enable users to refine queries in an intuitive and dynamic fashion while
81 exploring relationships in the data. Researchers can alter results along multiple search dimensions—
82 for subject, sample, and genomic criteria — without rerunning their query — by refining the
83 visualizations to samples with specific variants, or combinations of ranges and categories in
84 phenotypic and demographic dimensions. As the user changes the extent of any one category, that
85 change is projected to all other displayed data.

86 Scientists can use their intuition to generate hypotheses, then quickly look for initial
87 confirmation, and readily refine the direction of their search to pursue a suspected causal effect. The
88 resulting synthetic cohorts resolve to a set of individual subjects and samples, and the contents of
89 that set can be extracted for further analysis. A frequent current use case comes from our
90 collaboration with the Target ALS Multicenter Human Postmortem Tissue Core [11]: ALS
91 researchers can identify decedent biosamples of specific interest to their research using
92 MetroNome's data exploration capabilities, and then work with Target ALS Core directors to
93 rapidly obtain the specific blinded sample sets culled from multiple academic centers necessary for
94 rigorous follow-up experiments.

95

96 Comparison of cohorts

97 We provide the ability to view two cohorts side-by-side, to allow comparisons of traits that might
98 influence results and warrant further study. Visually comparing datasets is one way to determine
99 whether given cohorts are of interest, and whether specific dynamic filters better isolate features of
100 interest. Cohort-normalized values, such as gene expression z-scores, are recalculated on the fly,
101 based on the user's search parameters, to accurately represent critical differences between groups,
102 e.g., cases and controls.

103

104 Example use case

105 To illustrate the utility of MetroNome, we can use examples from ALS research: Figure 1 is a
106 comparison between patients with (left) and without (right) C9ORF72 repeat expansions, showing

107 gene expression patterns for the gene FIG4. The neuro axis diagrams present clear differences
 108 between the two groups in the primary motor cortex and, for the cohort with repeat expansions
 109 (left), the motor cortex vs. occipital cortex, an uninvolved region. In the relationship diagrams, this
 110 cohort is also marked by shorter duration of the disease.

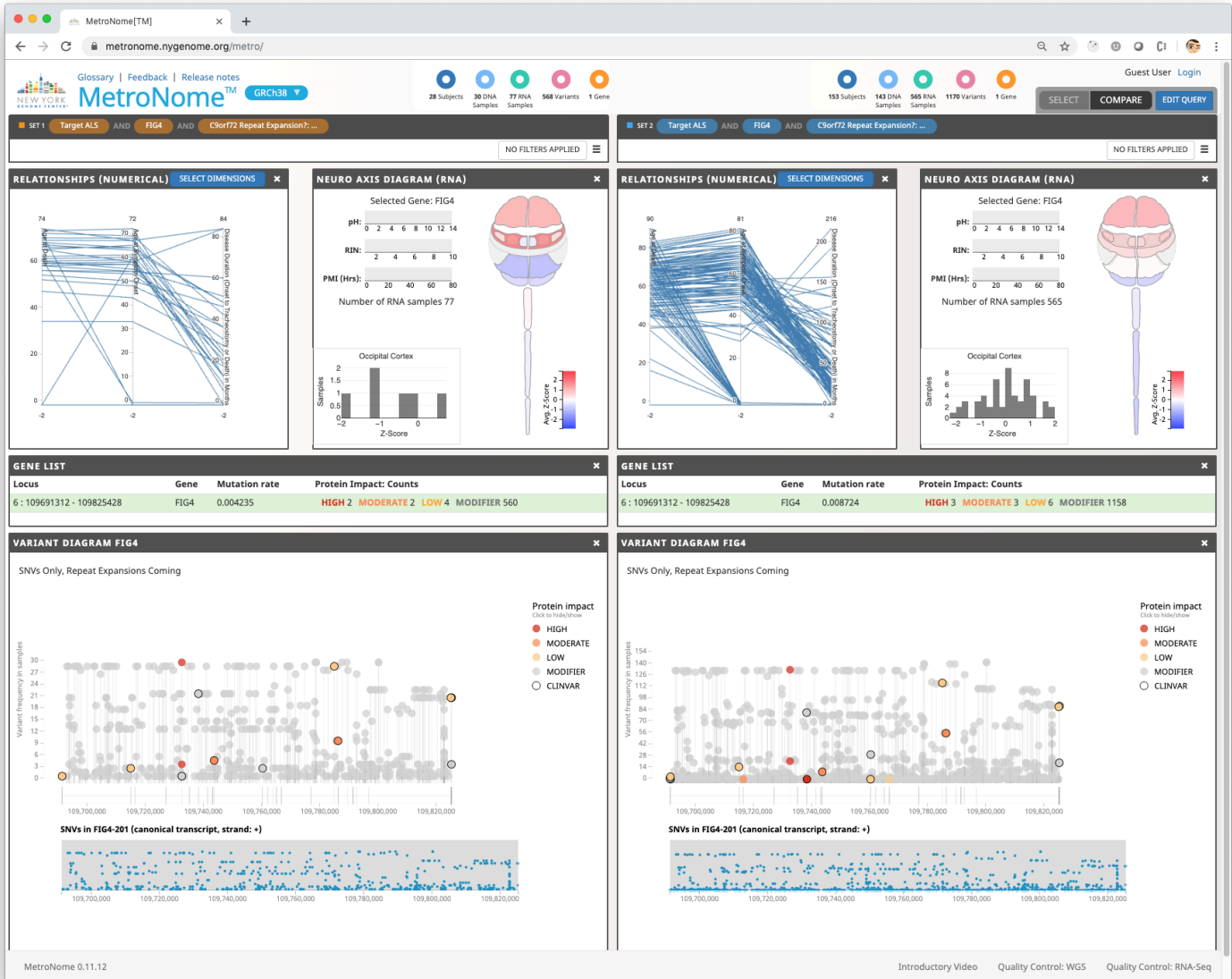


Figure 1: Comparison of cohorts with *C9ORF72* repeat expansions (left) and without (right), showing gene expression patterns for the gene FIG4 in an anatomogram (top), and variants (bottom).

111
 112 Figure 2 shows a query for samples with variants in the gene PFN1, requested by a researcher who
 113 had seen white matter abnormalities in mouse spinal cords. The neuro axis diagram clearly shows
 114 higher PFN1 expression in the spinal cord, and particularly in the thoracic spinal cord, which is
 115 interesting because the thoracic cord has a higher proportion of white matter compared to cervical
 116 and lumbar. In addition, the RNA heatmap indicates that there are a couple of specific samples with

117 very high cortical PFN_I expression. These specific decedent tissues and slides can be selected for
118 further benchtop experiments, and correlation with phenotypic information.
119



Figure 2: a query for samples with variants in the gene PFN_I. The neuro axis diagram (top) clearly shows higher PFN_I expression in the spinal cord, and particularly in the thoracic spinal cord. The RNA heatmap (bottom), besides showing generally higher expression in the spinal cord samples, reveals a couple of samples with very high cortical levels.

120

121

122 These types of searches can be performed instantly and without the need for prior bioinformatics
123 training. Lists of subjects, samples, and variants can be downloaded for further analysis and used to
124 identify and request tissues/biofluids/slides meeting specific criteria (such as gene expression
125 patterns, QC measures, or specific variants) from participating biorepositories for further benchtop
126 experimentation.

127 Infrastructure to support cross-study dynamic visualization

128 In addition to the visualization tools themselves, a great deal of data infrastructure is needed to
129 realize the outlined goals. These infrastructure tools center on integrating data with available
130 standards and with other data across multiple studies. We outline here three areas we have found
131 necessary to address:

- 132 • **Data harmonization:** We harmonize data where possible, which is currently a manual process.
133 Fields with different names in different datasets must be changed to the same standard names in
134 each. The values of those fields must be converted to the same vocabulary, the same numeric
135 ranges, and the same units. Values that are in reference to a specified or assumed range, such as
136 values from some laboratory tests that have varying ranges by instrument used, must be recorded
137 along with the reference range. Missing values must be handled, whenever possible, without
138 discarding the entire record.
- 139 • **Use of metadata and provenance:** the data needed for analyzing results for a single study are
140 often insufficient for integrating that data with information from other sources. Where the
141 metadata and provenance data exist, we use such information to more accurately present
142 combined data. We flag uncertainties to minimize misleading results.
- 143 • **Reference data:** to enable as much data integration as possible, we maintain significant types
144 of reference data, including target sets for standard exome kits. ClinVar [clinvar.org]
145 annotations are used to filter for pathogenic variants. Ensembl transcript-level and protein
146 domain annotations provide information on protein-coding impact and high-impact variants,
147 when disease significance is still unknown. 1,000 Genomes and TCGA somatic data are
148 maintained, primarily as sources of additional data for sparse datasets. These reference data are
149 used for comparisons, for interpretation of metadata, and for annotation of the synthesized
150 cohorts generated in MetroNome. Note that full data integration to enable further analysis is a
151 much more extensive problem that we have yet to address. Our work here is focused on enabling
152 integrated and comparative visualizations.

153 Privacy and security

154 Privacy and security are major concerns when we are supporting limited-access datasets.
155 Authorization to access a particular dataset is determined by the owner of that dataset. The NYGC
156 Data Privacy Committee must review the owner's approval before access is granted within
157 MetroNome. While this is a somewhat burdensome process, we feel it necessary to ensure that we
158 can host private data without risk of unintentional disclosure. If a user is approved for access, they
159 can grant members of their lab further access without review. This last case happens most often
160 when a user is part of a consortium, and the official owner of the data is their institution. The one
161 exception to that rule is that we allow users to upload any dataset to which they already have access,
162 limited to personal access only, and combine it with other MetroNome datasets they can access.

163 The MetroNome backend adds a clause to every database query to restrict the query to the set of
164 samples to which the user has access. Currently, users can access public data without logging in. If
165 they do so, the backend defaults to the “public” user, with associated access rights, and limits all
166 queries to public data only. Thus, no query or request can bypass the front end and avoid the
167 privacy checks.

168 The MetroNome front-end runs in an isolated subnet, accessible from outside the firewall. The
169 database and middle tier run behind the firewall, and all data thus resides internally. There is a
170 single connection through that firewall, restricted to a single machine address. Authentication is
171 performed for each connection established.

172 Technical implementation

173 Architecture and technologies used

174 MetroNome is built on a column-oriented database, Vertica, that holds all data (variants, gene
175 expression, phenotypic, demographic, and sample-related data). An application server, written in
176 Java, processes requests to the database and prepares results for display in the UI. The frontend is
177 written in JavaScript, using the React framework and D3 visualization library, and is hosted in a
178 TomCat web server.

179 System requirements

180 The system is designed to operate on Linux virtual machine nodes running CentOS 7 and can be
181 scaled to meet demands. The current database, Vertica, requires a cluster of dedicated nodes with
182 matching specifications for best operation; that is, each node of the cluster should be similar in
183 CPU, clock speed, number of cores, memory, and operating system version.

184 Source code and availability

185 The source code is currently being converted to open source and will be made available on
186 GitHub: <https://github.com/nygenome/metronome>.

187 The MetroNome installation hosted at the New York Genome Center is publicly available at
188 <https://metronome.nygenome.org>

189 Current use: Target ALS

190 The Target ALS Resource Cores were conceived to accelerate ALS therapy development by
191 providing the necessary highly curated biosamples and data resources broadly and rapidly to the
192 entire ALS research community. Given the numerous failures in translating laboratory results into
193 clinically effective therapies, a crucial aim was to address the substantial unmet need for high quality
194 patient-derived biosamples – such as brain, spinal cord, and muscle tissue samples from patients who
195 died from ALS and controls, and biofluids and stem-cell lines collected during disease progression.

196 We perform centralized Whole Genome Sequencing (WGS), and RNA-Seq for multiple central
197 nervous system regions at the New York Genome Center on every autopsy performed at one of the

198 academic centers comprising the federated Target ALS Postmortem Tissue Core. After passing QC,
199 the clinically annotated genomic and transcriptomic data is ingested into MetroNome and remains
200 linked to the tissue samples and de-identified metadata via Global Unique Identifiers (GUIDs). The
201 WGS and RNA-Seq raw data files (in multiple formats) are also made immediately available without
202 embargo or IP concerns – via an online form and established data transfer workflow.

203 MetroNome enables researchers with very little background in genetic analysis to access the data
204 set in a meaningful way, with relevance to their personal research interests. When researchers are
205 interested in a specific pathway or target, they can use MetroNome to search for cases with specific
206 mutations or variants, explore new hypotheses by comparing different cohorts defined by anatomical
207 regions or subject groups, and then directly request tissue samples from those cases. For virtually
208 every research project utilizing human biosamples, MetroNome can be used to refine slide and
209 sample sets, and direct further analysis. As examples, MetroNome
210 can be used to

- 211 • identify relevant tissue samples or slides for further benchtop experiments;
- 212 • find variants or RNA expression changes in specific targets;
- 213 • provide “clean controls” that do not possess mutations or unknown variants;
- 214 • compare spatial expression patterns with published imaging biomarker data meant to
215 quantify relevant pathways;
- 216 • examine whether gene expression patterns are consistent with activation of pathways
217 modulated by potential new drug candidates;
- 218 • identify whether specific patient subgroups display gene signatures that might inform
219 patient selection for clinical trials (Figure 3);
- 220 • segregate patients based on spatial gene expression patterns and correlate with fast/slow
221 progressors, site of onset, or specific neuropathological metadata;
- 222 • design further collaborative analysis of the genetic raw data and samples, such as whether
223 subject groups with distinct genetic patterns might correlate with biomarker profiles in
224 fluids or peripheral tissues.

225
226 The MetroNome visual data exploration platform has proven critical to the continued success,
227 expansion, and evolution of the Target ALS Postmortem Core and associated efforts. It has
228 supported over 100 different academic and industry labs, facilitating more than 135 different ALS
229 research projects in 16 countries across 4 continents. This often includes multiple different projects
230 in each lab. MetroNome has become part of a scientific ecosystem that includes clinics, research
231 labs, and industry.

232 General use

233 MetroNome is unbiased towards specific disease areas and can accommodate genomic and
234 phenotypic data from any study. Figure 3 shows an example from an esophageal cancer study,
235 indicating presence of single-nucleotide variants, copy number variations, and structural variations
236 for a matrix of genes and samples.

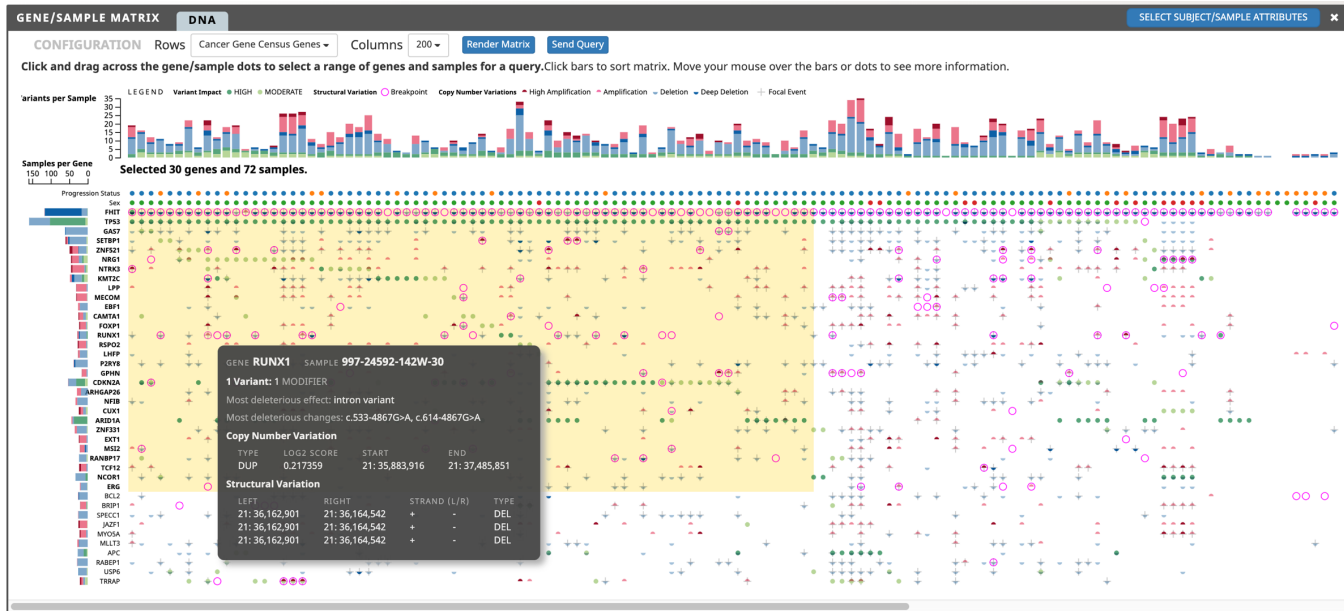


Figure 3: Sortable gene/sample matrix for genes from the Cancer Gene Census, shown in rows; columns represent samples. The grey callout shows details for one gene in one sample. Single-nucleotide variants are identified as high or moderate impact. Copy number changes and structural variants are identified by special glyphs and colors. Sorted by default to show genes and samples with the most impactful variants, this matrix can be used to select samples and genes (yellow area), e.g., to create a synthetic cohort prioritizing highly mutated samples.

237

238 Future work

239 New releases will include some additional critical features:

- 240 • Download images from our visualizations, along with the metadata about the cohort
241 being used and the filters applied.
- 242 • Automate the process for users to upload their own datasets, visible only to them.
- 243 • Improve reproducibility: users will be able to save queries, to rerun them in a later
244 session, or share them with other users. The ability to save *results*, and access them later,
245 should underlying data change in the interim, is planned. Finally, tracking the steps a
246 user executes in a session, displaying that history and allowing a user to return to a
247 previous state is a feature we believe to be very useful in this context.
- 248 • Add support and visualizations for new data types, such as repeat expansions and splice
249 junctions. We expect the types of data to expand continually.
- 250 • Automate harmonization using HPO terms [15].
- 251 • Pedigree relationships, including flagging filtered de novo or recessive homozygous
252 variants in the probands. Our structure allows subjects to be considered probands in one
253 study and relatives in others.

- 254 • To aid evaluation of results, we want to integrate statistical analysis tools. We may link
255 out to R, to BioConductor tools [www.bioconductor.org], or to tools such as DeepSea
256 [16].
- 257 • APIs to allow MetroNome to exchange both data and compute with other repositories.
258 These APIs are essential to users who wish to use the MetroNome resources with their
259 automated analyses, rather than through our visualization interface.
- 260 • Develop user interfaces for longitudinal data

261 Summary

262 Synthesizing cohorts by integrating data from multiple studies presents numerous challenges.
263 Providing this functionality as part of an interactive phenotype-genotype visualization platform
264 enables data integration as a fundamental part of the platform. Not only does this approach enhance
265 integrated multi-modal analysis, it provides a framework that reduces the work each researcher must
266 perform to obtain a clean cohort that meets their research needs. Using this visual data integration
267 platform to generate and explore hypotheses is a further important contribution, with the potential
268 to accelerate the work of scientists anywhere by eliminating the bioinformatics bottleneck during
269 genesis of ideas. Researchers can then follow up *only the best leads* with their computational colleagues
270 for thorough analysis.

271 Acknowledgements

272 The authors would like to thank the Target ALS Multicenter Postmortem Tissue Core for
273 making their sample data publicly available, and for valuable input during development of
274 functionality in MetroNome; Kanika Arora and Minita Shah for evaluating MetroNome's utility for
275 cancer studies and for providing extensive feedback during development of the gene/sample matrix;
276 Dorian Leary, Dimitrije Jevremović, Joseph Mulvaney, and Sylvestre Gug for their software
277 engineering; thanks to Phaedra Agius, Michael Zody, Simon Tavaré, and Tom Maniatis for
278 reviewing the manuscript.

279 References

- 280 1. Kahn SD. On the Future of Genomic Data. *Science* 2011; 331 (6018), 728–9
281 <https://doi.org/10.1126/science.1197891> PMID: 21311016
- 282 2. International Society for Biocuration (2018) Biocuration: Distilling data into knowledge. *PLoS*
283 *Biol* 16(4): e2002846. <https://doi.org/10.1371/journal.pbio.2002846>
- 284 3. McMurry JA, Koehler S, Washington NL, Balhoff JP, Borromeo C, Brush M, et al. Navigating
285 the Phenotype Frontier: The Monarch Initiative. *Genetics*. 2016 Aug; 203(4):1491–5.
286 <https://doi.org/10.1534/genetics.116.188870> PMID: 27516611.

- 287 4. Alyass et al. From big data analysis to personalized medicine for all: challenges and
288 opportunities. *BMC Medical Genomics* (2015) 8:33. <https://doi.org/10.1186/s12920-015-0108-y>
- 289 5. B. Shneiderman (1994) [Dynamic queries for visual information seeking](#). *IEEE software* 11 (6), 70-
290 77
- 291 6. Ritchie et al. Methods of integrating data to uncover genotype-phenotype interactions. *Nature*
292 *Reviews Genetics* | AOP, 13 January 2015; <https://doi.org/10.1038/nrg3868>
- 293 7. E.M. Ramos et al. (2013) Phenotype-Genotype Integrator (PheGenI): synthesizing genome-wide
294 association study (GWAS) data with existing genomic resources. *European Journal of Human*
295 *Genetics*; <https://doi.org/10.1038/ejhg.2013.96>
- 296 8. K.M. Wong et al. The dbGaP data browser: a new tool for browsing dbGaP controlled-access
297 genomic data, *Nucleic Acids Research*, Volume 45, Issue D1, 4 January 2017, Pages D819-
298 D826, <https://doi.org/10.1093/nar/gkwt139>
- 299 9. Pablo Cingolani, Adrian Platts, Le Lily Wang, Melissa Coon, Tung Nguyen, Luan Wang, Susan
300 J. Land, Douglas M. Ruden and Xiangyi Lu, A program for annotation and predicting the effects
301 of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster*
302 strain w1118; iso-2; iso-3, *Fly 6:2*, 1-13; April/May/June 2012
- 303 10. Landrum MJ, Lee JM, Benson M, Brown GR, Chao C, Chitipiralla S, Gu B, Hart J, Hoffman D,
304 Jang W, Karapetyan K, Katz K, Liu C, Maddipatla Z, Malheiro A, McDaniel K, Ovetsky M,
305 Riley G, Zhou G, Holmes JB, Kattman BL, Maglott DR. ClinVar: improving access to variant
306 interpretations and supporting evidence. *Nucleic Acids Res.* 2018 Jan 4.
- 307 11. L.W. Ostrow et al. [Target ALS Multicenter Human Postmortem Tissue Core](#). *ANNALS OF*
308 *NEUROLOGY* 76, S65-S65
- 309 12. Abecasis, Gonçalo R., A global reference for human genetic variation, *Nature*, 526: 68,
310 <https://doi.org/10.1038/nature15393>, 9/30/ 2015.
- 311 13. [The Cancer Genome Atlas Research Network](#), [John N Weinstein](#), [Eric A Collisson](#), [Gordon B](#)
312 [Mills](#), [Kenna R Mills Shaw](#), [Brad A Ozenberger](#), [Kyle Ellrott](#), [Ilya Shmulevich](#), [Chris](#)
313 [Sander](#) & [Joshua M Stuart](#), The Cancer Genome Atlas Pan-Cancer Analysis Project, *Nature*
314 *Genetics* volume 45, pages 1113-1120 (2013)
- 315 14. Jianjiong Gao, Bülent Arman Aksoy, Ugur Dogrusoz., Gideon Dresdner., Benjamin Gross.,
316 [Chris Sander](#), Nikolaus Schultz, Integrative Analysis of Complex Cancer Genomics and
317 Clinical Profiles Using the cBioPortal, *Sci. Signal.* 02 Apr 2013, Vol. 6, Issue 269, pp. pii, DOI:
318 10.1126/scisignal.2004088
- 319 15. [Robinson PN](#), [Mundlos S.](#), The human phenotype ontology., *Clin Genet.* 2010 Jun;77(6):525-34.
320 doi: 10.1111/j.1399-0004.2010.01436.x. Epub 2010 Feb 11, PMID: 20412080
- 321 16. Zhou, [Jian](#) & [Olga G Troyanskaya](#), Predicting effects of noncoding variants with deep learning-
322 based sequence model, *Nature Methods* volume 12, pages 931-934, 24 August 2015