

Supplementary Information: The Canadian VirusSeq Data Portal & Duotang: open resources for SARS-CoV-2 viral sequences and genomic epidemiology

1. Supplementary Methods

1.1 DNASTack Viral AI network for genomic variant surveillance

Viral AI is the world's first federated network for genomic variant surveillance, developed by DNASTack in response to the COVID-19 pandemic to support discovery and access to SARS-CoV-2 data. DNASTack partnered with the CanCOGeN - VirusSeq project to make the VirusSeq data available through Viral AI and to support lineage assignment for tracking variants.

Viral AI introduces a new way to share and analyze genomics, clinical, administrative, and related data, facilitating insights about transmission, severity, diagnostics and vaccine escape. As an alternative to the centralized model, where data is uploaded to a single vendor-managed database, Viral AI adopts a federated architecture to connect, analyze, and share data without moving it. This model enables faster, more efficient, regulatory compliant, and regionally sovereign data management, enabling viral surveillance efforts to be more equitable, scalable, and sustainable (see figure S1).



Figure S1: Federation makes it possible to drive discoveries across distributed data without moving it.

Viral AI accelerates science by making data uniformly accessible through a user-friendly graphical interface and powerful programmatic interfaces, integrating data across different sources from around the world alongside VirusSeq, such as NCBI Sequence Read Archive (SRA) and European Center for Disease Prevention and Control, among others. Over one million viral sequences have been added with corresponding assemblies, variant calls, and lineage assignments, all harmonized through an open source bioinformatics pipeline.

Viral AI is powered by a software suite that is compliant with multiple GA4GH standards and facilitates responsible and interoperable genomic and biomedical data sharing including the [Data Connect](#), [Data Repository Service](#), [Service Registry](#), and [Service Info](#) standards.

33

34 **Publisher** is a data integration and sharing studio that enables data custodians to connect
35 any dataset, from any source, without moving it. Data custodians who contribute data retain
36 administrative control and have transparency into how it's used. DNASTack has connected a
37 number of open-source viral genomic data sets using Publisher alongside the VirusSeq
38 data.

39

40 **Explorer** is a federated data hub that makes it easier for researchers to find, access, and
41 analyze shared data. With Explorer, researchers can search and perform analyses across a
42 universe of connected datasets through a single user interface. The VirusSeq data is made
43 available in Explorer for researchers to discover, access, and analyze alongside the other
44 connected data sets.

45 **1.2 Lineage Assignment Pipeline**

46 An open-source bioinformatics pipeline was developed to run lineage assignment on the
47 SARS-CoV-2 genome assemblies obtained from VirusSeq. The resulting lineage
48 assignments, in combination with sample metadata and assemblies, are imported into Viral
49 AI where they are made available over GA4GH standard interfaces.

50

51 Assembled SARS-CoV-2 genomes are periodically retrieved from VirusSeq and lineages are
52 assigned using pangolin (Phylogenetic Assignment of Named Global Outbreak LINEages), a
53 tool developed to implement the Pango nomenclature for SARS-CoV-2 lineages. To ensure
54 that lineage assignments are as accurate as possible, the more accurate but slower UShER
55 mode of pangolin is used to assign lineage. Additionally, since pangolin nomenclature and
56 designations are continuously updated as new variants are sequenced and categorized,
57 both pangolin and its underlying databases are updated in sync with new releases. Upon
58 update to pangolin or its databases, all previously assigned lineages are re-assigned using
59 the most up-to-date databases.

60

61 In addition to assigning lineage, the pipeline also produces a single-line multifasta and, for
62 each assembly, the set of sites that differs from the SARS-CoV-2 reference genome. The
63 resulting metadata, variant sites, assemblies, and multifasta are processed through an
64 ingestion pipeline and connected to Viral AI where the data is made publicly available for
65 further analysis and interpretation. Following its ingestion into Viral AI, the lineage metadata
66 is retrieved and added to the VirusSeq Data Portal and remains crucial to the researchers
67 conducting variant surveillance.

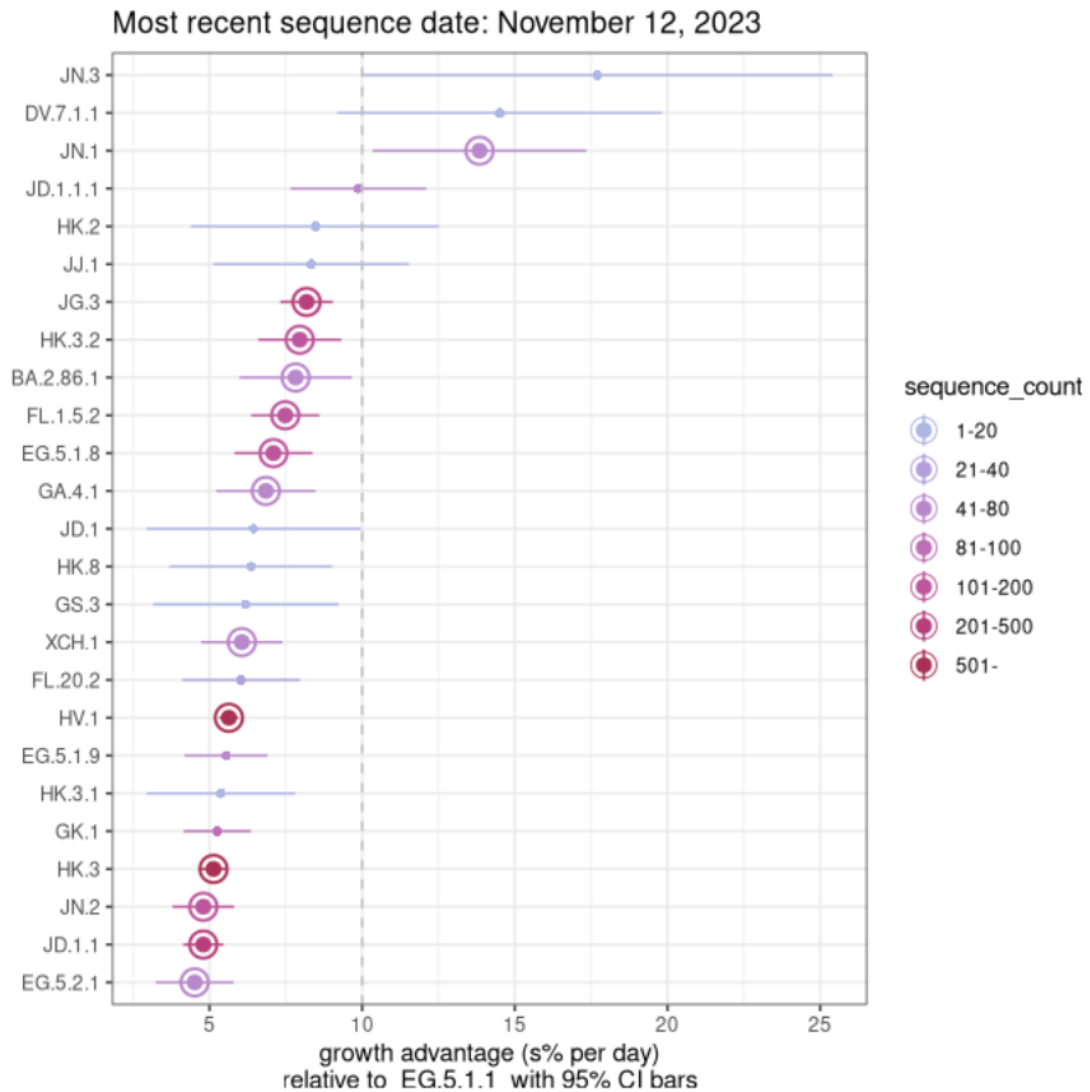
68 **2. Supplementary Results**

69 **2.1 List of Selected Contextual Data Fields Available on the Data Portal**

- 70 ● Study ID
- 71 ● Specimen Collector Sample ID
- 72 ● Sample Collected By
- 73 ● Sequence Submitted By
- 74 ● Submission Date
- 75 ● Sample Collection Date
- 76 ● Sample Collection Date Null Reason
- 77 ● Lineage Name
- 78 ● Lineage Analysis Software Name
- 79 ● Lineage Analysis Software Version
- 80 ● Lineage Analysis Software Data Version

- 81 ● Scorpio Call
- 82 ● Scorpio Version
- 83 ● Geo_loc_name (Country)
- 84 ● Geo_loc_name (State/province/territory)
- 85 ● Organism
- 86 ● Isolate
- 87 ● Fasta Header Name
- 88 ● Purpose Of Sampling
- 89 ● Purpose Of Sampling Details
- 90 ● Anatomical Material
- 91 ● Anatomical Part
- 92 ● Body Product
- 93 ● Environmental Material
- 94 ● Environmental Site
- 95 ● Collection Device
- 96 ● Collection Method
- 97 ● Host (Scientific Name)
- 98 ● Host Disease
- 99 ● Host Age
- 100 ● Host Age Null Reason
- 101 ● Host Age Unit
- 102 ● Host Age Bin
- 103 ● Host Gender
- 104 ● Purpose Of Sequencing
- 105 ● Purpose Of Sequencing Details
- 106 ● Sequencing Instrument
- 107 ● Sequencing Protocol
- 108 ● Raw Sequence Data Processing Method
- 109 ● Dehosting Method
- 110 ● Consensus Sequence Software Name
- 111 ● Consensus Sequence Software Version
- 112 ● Breadth Of Coverage Value
- 113 ● Depth Of Coverage Value
- 114 ● Reference Genome Accession
- 115 ● Bioinformatics Protocol
- 116 ● Gene Name
- 117 ● Diagnostic Pcr Ct Value
- 118 ● Diagnostic Pcr Ct Value Null Reason
- 119 *For a complete list of Data Portal policies and available contextual data, view*
- 120 <https://virusseq-dataportal.ca/policies>

Plot single lineages in Canada *

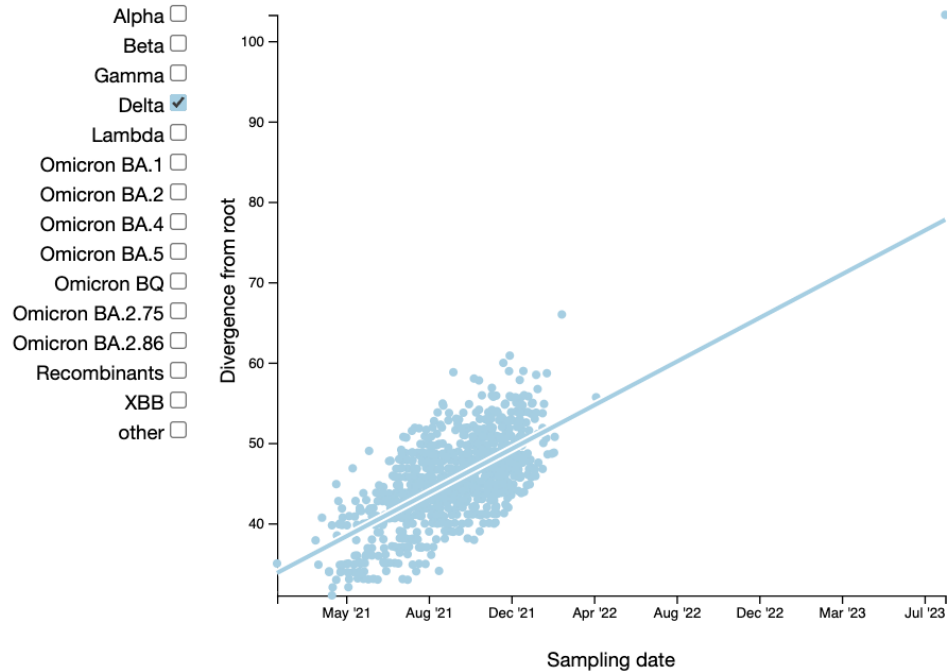


121
122
123
124
125
126
127
128

Figure S2: The Fastest Growing Lineages Plot allows identification of variants with likely true selective advantage vs. those that rise in frequency by chance. Selective advantage (with 95% credibility intervals) is calculated for each variant in each province. Those variants that have a positive growth advantage in more than one province are denoted with a circle around the plotted point in the overall Canadian plot. This plot thus shows the strength of evidence for selection both over time (reflected in smaller credibility intervals) and over jurisdictions.

Root-to-tip analyses

The slope of root-to-tip plots over time provide an estimate of the substitution rate. A lineage with a steeper positive slope than average for SARS-CoV-2 is accumulating mutations at a faster pace, while a lineage that exhibits a jump up (a shift in intercept but not slope) has accumulated more than expected numbers of mutations in a transient period of time (similar to what we saw with Alpha when it first appeared in the UK).



129
130
131
132
133
134

Figure S3: Root-to-tip analyses allow rapid identification of novel appearances of variants that were previously dominant. By examining a root-to-tip plot, Duotang users can quickly identify a sample that possesses more mutations than would be expected given the time since the lineage emerged, or samples that are collected much later than would be expected given the timing of a wave caused by a specific variant.

135 3. Consortium and Network Author Information

136 Canadian Public Health Laboratory Network (CPHLN) members and staffs having 137 contributed data to the portal

138 1. Alberta

139 *Genomes with prefix ABPHL*

140 Bu J, Croxen M, Deo A, Dieu P, Dong X, Ferrato C, Gavriliuc S, George R, Getachew F, Gill,
141 K, Ie N, Khadka R, Khan F, Koleva P, Lee L, Li V, Lindsay A, Lloyd C, Lynch T, Ma R,
142 McCullough, E, Mohon A, Murphy S, Obasuyi O, Pabbaraju K, Presbitero A, Rotich S,
143 Shokoples S, Thayer J, Tipples G, Trevor H, Whitehouse M, Wong A, Yu C, Zelyas N

144 *Genomes in collaboration with University of Calgary (prefix AB-NNNNNN)*

145 Gordon P, Lam LG, Pabbaraju K, Wong A, Ma R, Li V, Melin A, Tipples G, Berenger B,
146 Zelyas N, Kellner J, Bernier F, Chui L, Croxen M

147 **2. British Columbia**

148 *BCCDC Public Health Laboratory*

149 Natalie Prystajeky, Linda Hoang, John R. Tyson, Dan Fornika, Shannon Russell, Kim
150 MacDonald, Kimia Kamelian, Ana Pacagnella, Corrinne Ng, Loretta Janz, Richard Harrigan,
151 Robert Azana, Tara Newman, Jessica Caleta, Sherrie Wang, Janet Fung, Mel Krajden

152 **3. Manitoba**

153 *Cadham Provincial Laboratory collected specimens sequenced at Canada's National
154 Microbiology Lab (NML)***

155 Paul Van Caesele, Jared Bullard, David Alexander, Kerry Dust

156 *Cadham Provincial Laboratory sequenced specimens*

157 David Alexander, Lori Johnson, Janna Holowick, Joanne Sanders, Adam Hedley, Kerry Dust,
158 Ayo Bolaji, Brooke Cistarelli, Emma Rempel, Paul van Caesele, Jared Bullard

159 *Dynacare sequenced specimens*

160 Hilary Racher, Melissa Desaulnier, Tintu Abraham, Hongbin Li (Impact Genetics, Brampton
161 Ontario)

162 *NML***

163 See below for list of NML personnel

164 **4. New Brunswick**

165 *Centre Hospitalier Universitaire Georges L. Dumont*

166 Beauregard AP., Lyons P., Chacko S., Shaw W., Lacroix J., Allain E., Crapoulet N., Garceau
167 R., Desnoyers G.

168 *NML***

169 See below for list of NML personnel

170 **5. Newfoundland and Labrador**

171 *Newfoundland and Labrador Health Services*

172 *(previously known as "Newfoundland and Labrador - Eastern Health")*

173 Robert Needle, Yang Yu, Laura Gilbert, George Zahariadis, Geoffrey Woodland, Chris
174 Corkum, Kerri Smith, Phillip Andrews, Matthew Gilmour

175 *NML***

176 See below for list of NML personnel

177 **6. Nova Scotia**

178 *QEII Health Sciences Centre**

179 Todd Hachette, Jason LeBlanc, Janice Pettipas, Dan Gaston, Greg McCracken
180 (*data tagged post September 14th should include Allana Loder)

181 *NML***

182 See below for list of NML personnel

183 **7. Ontario**

184 *Public Health Ontario Laboratory*

185 Jacob Afelskie, Vanessa G Allen, Rebecca Azzaro, Doonia Bajovic, Philip Banh, Ilse
186 Belgrave, Tom Braukmann, Ashley Carandang, Yao Chen, Claudia Chu, Shawn Clark, Kirby
187 Cronin, Richard de Borja, Rachelle DiTullio, Carla Duncan, Hadi El Roz, Alireza Eshaghi,
188 Nahuel Fittipaldi, Christine Frantz, Dhiraj Gaglani, Nicole Graham, Jonathan B Gubbay,
189 Jennifer L Guthrie, Lawrence Heisler, Daniel Heydari, Mark Horsman, Hadia Hussain, Jason
190 Iraheta, Grace Jeong, Esha Joshi, Sushma Kavikondala, Lisa Kim, Surendra Kumar, Michael
191 Laszloffy, Aimin Li, Michael C.Y. Li, Alex Marchand-Austin, Maria Mariscal, Dean Maxwell,
192 Lisa McTaggart, Fatima Merza, Anupam Mittal, Naadia Mohammed, Esther Nagai, Sandeep
193 Nagra, Shiva Nassori, Paul Nelson, Rima Palencia, John Palmer, Samir N Patel, Stephen
194 Perusini, Nataliya Potapova, Anna Puzinovi, Zarah Rajaei, Christina Rampertab, Himeshi
195 Samarsinghe, Candice Schreiber, Christine Seah, Fatemeh Shaeri, Kapil Shaeri, Kapil Shah,
196 Narisha Shakuralli, Natasha Singh, Karthikeyan Sivaraman, Brenda Stanghini, Ashleigh
197 Sullivan, Vincent Su Bin Cha, Yogi Sundaravadanam, Sarah Teatero, Semra Tibebu, Nobish
198 Varghese, Andre Villegas, Jesse Wang, Matthew Watson, Sichong Xu, Xiao Xu, Kent Young,
199 Sophie Yu, Farhan Yusuf, Sandra Zittermann

200 *Ontario Institute for Cancer Research*

201 Jared T. Simpson, Richard de Borja, Paul Krzyzanowski, Bernard Lam, Lawrence Heisler,
202 Michael Laszloffy, Yogi Sundaravadanam, Ilinca Lungu, Lubaina Kothari, Cassandra
203 Bergwerff, Jeremy Johns, Felicia Vincelli, Philip Zuzarte

204 *McMaster University*

205 Hooman Derakhshani, Sheridan J.C. Baker, Emily M. Panousis, Ahmed N. Draia, Jalees A.
206 Nasir, Michael G. Surette, Andrew G. McArthur

207 **8. Prince Edward Island**

208 *Queen Elizabeth Hospital*

209 Xiaofeng Ding, Vanessa Arseneau, Kari-Lyn Young

210 *NML***

211 See below for list of NML personnel

212 **9. Québec**

213 *Laboratoire de Santé Publique du Québec, McGill Génome Sciences Centre; CoVSeq*
214 *Consortium*

215 Sandrine Moreira, Jiannis Ragoussis, Guillaume Bourque, Éric Fournier, Aurélie Guilbault,
216 Benjamin Delisle, Dihya Baloul, Inès Levade, Sarah Reiling, Hector Galvez, Paul
217 Stretenowich, Alexandre Montpetit, Michel Roger, Judith Fafard

218 **10. Saskatchewan**

219 *Roy Romanow Provincial Laboratory*

220 Ryan McDonald, Keith MacKenzie, Kara Loos, Stefani Kary, Meredith Faires, Guruprasad
221 Janga, Rachel DePaulo, Laura Klassen, Alanna Senecal, Amanda Lang, Jessica Minion,
222 Roy Romanow Provincial Laboratory - Molecular Diagnostics

223 *NML***

224 See below for list of NML personnel

- 225 **11. Canada's National Microbiology Laboratory (NML)****
- 226 Anna Majer, Shari Tyson, Grace Seo, Philip Mabon, Elsie Grudeski, Rhiannon Huzarewich,
227 Russell Mandes, Anneliese Landgraff, Jennifer Tanner, Natalie Knox, Morag Graham, Gary
228 Van Domselaar, Nathalie Bastien, Ruimin Gao, Cody Buchanan, Jasmine Frost, Ameet
229 Bharaj, Cole Slater, Nikki Toledo, Laura Hart, Yan Li, Timothy Booth, Catherine Yoshida,
230 Genevieve Labbe, Adina Bujold, Kara Loos, Jennifer Beirnes, Michael Przybytkowski,
231 Patrick Bastedo, Debra Sorensen, Andrea Tyler, Ana Duggan, Darian Hole, Madison Chapel,
232 Kristen Biggar, Emily Haidl, Chanchal Yadav, Jeff Tuff, Connor Chato, Katherine Eaton,
233 Kirsten Palmier, Molly Pratt, Amber Papineau, Adrian Zetner, Carmen Lia Murall
- 234 Genomics Core Facility at NML
235 Robotics Support Laboratory at NML
236 **Academic, Health Network and Research Institutions staffs having contributed data to the**
237 **Portal**
- 238 *Unity Health Toronto*
239 Ramzi Fattouh, Larissa M. Matukas, Yan Chen, Mark Downing, Trina Otterman, Karel
240 Boissinot, Le Luu
- 241 *University Health Network/Mount Sinai Hospital Department of Microbiology*
242 Marie-Ming Aynaud, Javier Hernandez, Seda Barutcu, Kin Chan, Jessica Bourke, Marc
243 Mazzulli, Tony Mazzulli, Laurence Pelletier, Jeff Wrana, Aimee Paterson, Angel Liu, Allison
244 McGeer
- 245 *Kingston Health Sciences Centre and Queen's University*
246 Prameet M. Sheth, Calvin Sjaarda, Robert Colautti, Katya Douchant
- 247 *University of British Columbia*
248 John R. Tyson, Gabrielle Jayme, Karen Jones, Terrance P. Snutch
- 249 *Toronto Invasive Bacterial Diseases Network; Sunnybrook Health Sciences*
250 Allison McGeer, Patryk Aftanas, Angel Li, Kuganya Nirmalarajah, Emily Panousis, Ahmed
251 Draia, Jalees Nasir, David Richardson, Michael Surette, Samira Mubareka, Andrew G.
252 McArthur
- 253 *Eastern Ontario Regional Laboratory Association*
254 Leanne Mortimer, Hooman Derakhshani, Emily Panousis, Ahmed Draia, Jalees Nasir, Robert
255 Slinger, Andrew G. McArthur
- 256 *Western University, CoVizu^{rs} Dev Team*
257 Art Poon, Gopi Gugan, Bonnie Lu, Roux-Cil Ferreira, Molly Liu, Laura Muñoz Baena, Kaitlyn
258 Wade, Navaneeth Mohan, Sandeep Thokala, Abayomi Olabode
- 259 **CanCOGeN VirusSeq committees and working groups' members**
- 260 *CanCOGeN VirusSeq Implementation Committee*
261 Terrance Snutch, Fiona Brinkman, Marceline Côté, William Hsiao, Gordon Jolly, Yann Joly,
262 Sharmistha Mishra, Sandrine Moreira, Samira Mubareka, Jared Simpson, Megan
263 Smallwood, Gary Van Domselaar
- 264 *CanCOGeN Capacity Building Working Group*
265 Gary Van Domselaar, Matthew Croxen, Natalie Knox, Celine Nadon, Jennifer Tanner

266 *CanCOGeN Data Analytics Working Group*
267 Gary Van Domselaar, Fiona Brinkman, Zohaib Anwar, Robert Beiko, Matieu Bourgey,
268 Guillaume Bourque, Richard de Borja, Ahmed Draia, Jun Duan, Marc Fiume, Dan Fornika,
269 Eric Fournier, Erin Gill, Paul Gordon, Emma Griffiths, Jose Hector Galvez Lopez, Darian
270 Hole, William Hsiao, Jeffrey Joy, Kimia Kamelian, Natalie Knox, Philip Mabon, Finlay
271 Maguire, Tom Matthews, Andrew McArthur, Samir Mechai, Sandrine Moreira, Art Poon,
272 Amos Raphenya, Claire Sevenhuysen, Jared Simpson, Jennifer Tanner, Lauren Tindale,
273 John Tyson, Geoff Winsor, Nolan Woods, Matthew Croxen, Carmen Lia Murall

274 *CanCOGeN Ethics and Governance Working Group*
275 Yann Joly, Fiona Brinkman, Erin Gill, William Hsiao, Hanshi Liu, Sandrine Moreira, Gary Van
276 Domselaar, Ma'n Zawati, Sarah Savić-Kallesø

277 *CanCOGeN Metadata Working Group*
278 William Hsiao, David Alexander, Zohaib Anwar, Nathalie Bastien, Tim Booth, Guillaume
279 Bourque, Fiona Brinkman, Hughes Charest, Caroline Colijn, Matthew Croxen, Guillaume
280 Desnoyers, Rejean Dion, Damion Dooley, Ana Duggan, Leah Dupasquier, Kerry Dust, Eleni
281 Galanis, Emma Garlock, Erin Gill, Gurinder Gopal, Tom Graefenhan, Morag Graham, Emma
282 Griffiths, Linda Hoang, Naveed Janjua, Jeffrey Joy, Kimia Kamelian, Lev Kearney, Natalie
283 Knox, Theodore Kuschak, Jason LeBlanc, Yan Li, Anna Majer, Adel Malek, Ryan McDonald,
284 David Moore, Celine Nadon, Samir Patel, Natalie Prystajacky, Anoosha Sehar, Claire
285 Sevenhuysen, Garrett Sorensen, Laura Steven, Lori Strudwick, Marsha Taylor, Shane
286 Thiessen, Gary Van Domselaar, Adrian Zetner

287 *CanCOGeN Research Collaborations Working Group*
288 Fiona Brinkman, Zohaib Anwar, Marceline Côté, Marc Fiume, Laura Gilbert, Erin Gill, Paul
289 Gordon, Yann Joly, Sandrine Moreira, Samira Mubareka, Natalie Prystajacky, Jennifer
290 Tanner, Gary Van Domselaar, Phot Zahariadis,

291 *CanCOGeN Sequencing Working Group*
292 Ioannis Ragoussis, Terrance Snutch, Patryk Aftanas, Matthew Croxen, Hooman
293 Derakhshani, Nahuel Fittipaldi, Morag Graham, Andrew McArthur, Sandrine Moreira, Samira
294 Mubareka, Natalie Prystajacky, Ioannis Ragoussis, Jared Simpson, Michael Surette, John
295 Tyson

296 *CanCOGeN Quality Control Working Group*
297 Jared Simpson, Mathieu Bourgey, Kodjovi Dodji Mlaga, Nahuel Fittipaldi, Jose Hector
298 Galvez Lopez, Natalie Knox, Genevieve Labbe, Pierre Lyons, Philip Mabon, Finlay Maguire,
299 Anna Majer, Andrew McArthur, Ryan McDonald, Sandrine Moreira, Natalie Prystajacky,
300 Karthikeyan Sivaraman, Kerri Smith, Terrance Snutch, Karthikeyan Sivaraman, Andrea Tyler,
301 John Tyson, Gary Van Domselaar, Matthew Croxen

302 *Canadian VirusSeq Data Portal (CVDP) Team*
303 Guillaume Bourque, Lincoln Stein, Christina Yung, Hanshi Liu, Yann Joly, Adrielle Houweling,
304 William Hsiao, Marc Fiume, David Bujold, Erin Gill, Fiona Brinkman, Nithu John, Rosita
305 Bajari, Linda Xiang, Alexandru Lepsa, Jaser Uddin, Justin Richardsson, Leonardo Rivera

306 Funding for the VirusSeq Data Portal is provided by The Canadian COVID Genomics
307 Network (CanCOGeN), and supported by Genome Canada and Innovation, Science and
308 Economic Development Canada (ISED)