## 1 TITLE

- 2 The Gastric Cancer Registry: A Genomic Translational Resource for Multidisciplinary
- 3 Research in Stomach Malignancies
- 4

## 5 AUTHORS

- 6 Alison F Almeda<sup>1</sup>, Sue Grimes<sup>1</sup>, HoJoon Lee<sup>1</sup>, Stephanie U Greer<sup>1</sup>, GiWon Shin<sup>1</sup>,
- 7 Madeline McNamara<sup>1</sup>, Anna C Hooker<sup>1</sup>, Maya Arce<sup>1</sup>, Matthew Kubit<sup>1</sup>, Marie Schauer<sup>1</sup>,
- 8 Paul Van Hummelen<sup>1</sup>, Cindy Ma<sup>1</sup>, Meredith Mills<sup>1</sup>, Robert J Huang<sup>2</sup>, Joo Ha Hwang<sup>2</sup>,
- 9 Manuel R Amieva<sup>3</sup>, Summer Han<sup>4</sup>, James M. Ford<sup>1</sup>, Hanlee P. Ji<sup>1\*</sup>
- <sup>1</sup>Division of Oncology, Department of Medicine, Stanford University School of Medicine,
- 11 Stanford, CA, 94305, United States
- 12 <sup>2</sup>Division of Gastroenterology and Hepatology, Department of Medicine, Stanford
- 13 University School of Medicine, Stanford, CA, 94305, United States
- <sup>3</sup>Division of Infectious Diseases, Department of Pediatrics, Stanford University School of
- 15 Medicine, Stanford, CA, 94305, United States
- <sup>4</sup>Department of Neurosurgery, Stanford University School of Medicine, Stanford, CA,
- 17 94305, United States
- 18

### 19 RUNNING TITLE

20 Developing the Gastric Cancer Registry and Genome Explorer

## 1 KEYWORDS

2 Gastric Cancer; Cancer Registry; Hereditary Cancer; Sequencing; Data Portal

3

## 4 FINANCIAL SUPPORT

- 5 The work was supported by the Gastric Cancer Foundation. Additional support to HPJ
- 6 came from the Research Scholar Grant, RSG-13-297-01-TBG from the American
- 7 Cancer Society, and the Clayville Foundation. RJH is supported by the National Cancer
- 8 Institute of the National Institutes of Health under Award Number K08CA252635.

9

## 10 CORRESPONDING AUTHOR

- 11 Hanlee P. Ji
- 12 Division of Oncology, Department of Medicine Stanford University School of Medicine
- 13 CCSR 2245, 269 Campus Drive
- 14 Stanford, CA 94305-5151
- 15 Email : genomics ji@stanford.edu
- 16 Fax : 650-736-1454
- 17 Phone : (650) 721-1503

18

## 19 CONFLICT OF INTEREST

20 The authors have no competing interests to declare.

21

## 22 MANUSCRIPT FEATURES

- 23 Abstract: 271 words
- 24 Manuscript: 3548 words
- 25 Tables: 2
- 26 Figures: 4

1 Supplementary files: 2

2

### 3 ABBREVIATIONS

- 4 GC Gastric cancer
- 5 GCR Gastric Cancer Registry
- 6 FFPE Formalin fixed paraffin embedded
- 7 RNA-seq RNA sequencing
- 8 HLA Human leukocyte antigen
- 9 TCGA The Cancer Genome Atlas
- 10 STAD Stomach adenocarcinoma
- 11 ESCA Esophageal carcinoma
- 12 StoP Stomach Cancer Pooling Project

#### 1 ABSTRACT

Background: Gastric cancer (GC) is a leading cause of global cancer morbidity and
mortality. Developing information systems which integrate clinical and genomic data
may accelerate discoveries to improve cancer prevention, detection, and treatment. To
support translational research in GC, we developed the GC Registry (GCR), a North
American repository of clinical and cancer genomics data.

Methods: GCR is a national registry with online self-enrollment. Entry criteria into the GCR included the following: (1) diagnosis of GC, (2) history of GC in a first- or seconddegree family member or (3) known pathogenic or likely pathogenic germline mutation in the gene *CDH1*. Participants provided demographic and clinical information through a detailed (412-item) online survey. A subset of participants provided specimens of saliva and tumor samples. These tumor samples underwent exome sequencing, whole genome sequencing and transcriptome sequencing.

**Results:** From 2011-2021, 567 individuals registered for the GCR and returned the 14 clinical guestionnaire. For this cohort 65% had a personal history of GC, 36% reported 15 16 a family history of GC and 14% had a germline CDH1 mutation. Eighty-nine GC 17 patients provided tumor tissue samples. For the initial pilot study, 41 tumors were sequenced using next generation sequencing. The data was analyzed for cancer 18 19 mutations, copy number variations, gene expression, microbiome presence, 20 neoantigens, immune infiltrates, and other features. We developed a searchable, web-21 based interface (the GCR Genome Explorer) to enable researchers access to these 22 datasets.

- 1 **Conclusions:** The GCR is a unique, North American GC registry which integrates both
- 2 clinical and genomic annotation.
- 3 Impact: Available for researchers through an open access, web-based explorer, we
- 4 hope the GCR Genome Explorer accelerates collaborative GC research across the
- 5 United States.

#### 1 INTRODUCTION

Gastric cancer (GC) is a leading cause of cancer morbidity and mortality worldwide.<sup>1</sup>
While incidence is lower in the United States, GC remains a major public health concern
with an estimated 116,000 prevalent cases nationwide in 2017.<sup>2</sup> GC is diagnosed at
generally advanced stages in the United States, where curative resection is often no
longer possible.<sup>2,3</sup> These data underscore the need for additional translational research
in GC etiology, prevention, early detection, and therapy.

8

9 Cancer registries are a valuable resource for collating clinical information. Some 10 registries also contain biological specimen repositories of both cancerous and non-11 cancerous tissue, allowing for somatic and germline genomic characterization through next generation sequencing.<sup>4,5</sup> Cancer registries that integrate clinical information with 12 tumor genomic features are particularly useful in translational research.<sup>6,7</sup> There 13 14 currently exist few registries focused on gastric cancer, particularly with patient data and samples derived from the United States. The availability of medical records, 15 16 epidemiological data, and biospecimens of tumor tissue allow researchers to analyze 17 genetic, environmental, and other differences that provide clues leading to risk factors. 18 This research is particularly relevant given the poor overall outcomes from GC in the 19 United States, and lack of established screening and surveillance programs for this 20 deadly cancer.

21

1 To address this knowledge gap, we established the Gastric Cancer Registry (GCR) in 2 2011. The goal of this project is to integrate granular patient clinical data (collected 3 through a detailed, 412-item online guestionnaire) with comprehensive genomic 4 characterization of tumor samples. This includes gene expression, somatic mutations, 5 copy number variation, human leukocyte antigen genotypes, neoantigens, and intra-6 tumoral heterogeneity details. To facilitate public access to these data, we created the 7 GCR Genome Explorer (https://gcregistry-explorer.stanford.edu/), a browser-based interactive tool which allows for querying of clinical and molecular annotation from the 8 9 GCR. In this manuscript, we describe the overall study design, methods for sample 10 collection and data generation, and characteristics of enrolled participants in the GCR 11 over a ten-year period (2011-2021). We also review features of the GCR Genome 12 Explorer and describe how this tool can be used by cancer researchers for translational 13 research.

14

#### 15 MATERIALS AND METHODS

#### 16 Study design

Approval was obtained from the Stanford University Institutional Review Board (IRB-20285). The study was performed in accordance with the ethical standards delineated in the 1964 Declaration of Helsinki and its later amendments. Eligible participants were over 18 years of age and fulfilled one or more of the following criteria: personal history of histologically proven GC, a family history of GC in a first- or second-degree relative, and/or a known pathogenic or likely pathogenic mutation in the gene *CDH1*, also known

1 as Cadherin-1. From 2011-2014, only individuals with a GC diagnosis were included in 2 the study. Family history and a known germline CDH1 mutation were added as 3 eligibility criteria in 2014. Biospecimens were collected from residents of the United 4 States. Recruitment and enrollment for the GCR was conducted through referral by 5 local clinicians, local and regional pamphlet distribution at conferences and 6 interest/advocacy group events, social media advertisement and a study website 7 (https://gcregistry.stanford.edu/). Registrants use the website to complete the 8 questionnaire.

9

10 Study data were collected and managed using REDCap electronic data capture tools hosted at Stanford University. REDCap is a secure, web-based software platform 11 12 designed to support data capture for research studies, providing 1) an intuitive interface 13 for validated data capture; 2) audit trails for tracking data manipulation and export 14 procedures; 3) automated export procedures for seamless data downloads to common statistical packages; and 4) procedures for data integration and interoperability with 15 external sources.<sup>8,9</sup> Moreover, REDCap allows for automated consent for participation 16 17 in research. Upon registration, participants were asked to complete a 412-item 18 questionnaire (**Supplementary File 1**). This intake survey is focused on information 19 relevant to genetics, lifestyle, environmental and personal risk factors related to GC.

20

The questionnaire contained questions regarding subjects' demographics, medical
history, familial cancer history, and lifestyle behaviors. For providing updates in

personal or family history, participants could provide additional information using a
unique return code. Participants had the option of authorizing the release of medical
records related to the diagnosis and treatment of the patient's GC. This information
included treatment summaries from any chemotherapy or radiation therapy, the
operative notes from any surgeries and any pathology reports relating to the diagnosis
of GC.

7

#### 8 Biospecimen collection

9 We obtained tissue specimens for a subset of participants. These samples included 10 archival tissue samples in the form of formalin fixed, paraffin embedded (FFPE) blocks, 11 scrolls and unstained slides. This cohort included individuals with GC who had a 12 surgical resection of the stomach or carriers of CDH-1 germline mutations who 13 underwent a screening endoscopic with biopsy or preventative gastrectomy. From 2018 14 onward, participants had the option of donating a sample of saliva. We used the Oragene DNA collection kit (DNA Genotek Inc., catalog no. OGR-500, OGR-600) to 15 16 collect these samples.

17

#### 18 Sample processing

Archival FFPE samples of gastric tumors were used for genomic studies. Genomic
DNA and total RNA were extracted from the tumor samples and the quality of nucleic
acids was determined through quantitation, genomic sizing, and fragmentation analysis.
We used an ultrasonicator to shear genomic DNA to a desired length of 300 base pairs.

1	The genomic DNA was transformed into sequencing libraries using the KAPA Hyper
2	Prep Kit for Illumina (Roche, catalog no. KK8502) with eight-base pair unique dual
3	adapters. An amount of genomic DNA library was set aside for exome enrichment
4	using xGen Lockdown Probes and reagents (Integrated DNA Technologies, catalog no.
5	1056114, 1075474). To prepare RNA sequencing libraries, the KAPA RNA HyperPrep
6	Kit (Roche, catalog no. KK8540) was used with modifications to accommodate low
7	quality samples. The genomic DNA, exome, and RNA libraries were pooled
8	respectively and sequenced on the NovaSeq 6000 system (Illumina, catalog no.
9	20012850). More details about the methods are described in <b>Supplementary File 2</b> .
10	
11	Next generation sequencing
12	The genomic DNA libraries underwent low-coverage 1-2X whole genome sequencing.
13	Exome libraries were sequenced more deeply at approximately 68X coverage to enable
14	somatic mutation calling. The RNA libraries were sequenced at high depth for an
15	average of 65 million reads per sample.

16

### 17 Bioinformatic analysis

Whole genome sequencing reads were analyzed for copy number variation. To identify somatic copy number changes for samples without a matched normal control, we used a normal reference genome data set as a comparison control.<sup>10</sup> Whole exome sequencing data was analyzed for single nucleotide variants and insertions/deletions of nucleotides in genomic DNA. RNA-seq data was investigated for multiple features,

including: (i) gene expression level, (ii) four-digit human leukocyte antigen (HLA)
genotypes, (iii) tumor immune infiltrate cell types, and (iv) the microbiome. Based on
whole exome sequencing data and RNA-seq, candidate cancer neoantigens were
identified. More detailed methods and bioinformatic pipeline specifics are in
Supplementary File 2.

6

#### 7 The GCR Genome Explorer

8 The GCR Genome Explorer contains clinical, genomic, and genetic data from the 9 registry. Within the GCR Genome Explorer, we also included other datasets derived 10 from the Cancer Genome Atlas (TCGA) stomach adenocarcinoma (STAD) and 11 esophageal carcinoma (ESCA) projects. The Variant Call Format files and copy number 12 variations for the TCGA cohorts were downloaded from the National Cancer Institute Genomic Data Commons.<sup>11</sup> Genomic and genetic TCGA data was processed using the 13 14 similar bioinformatics pipeline as GCR data (though with changes to accommodate for matched tumor-normal data), resulting in comparable bioinformatic analysis across all 15 16 samples.

17

The portal is a two-tier client-server application written using Ruby on Rails version
5.1.7 and ruby 2.4.9, with back-end database tables in MySQL 5.5.62 and deployed
using Passenger and Apache2. The application server has 64GB RAM and 32
processors running Ubuntu 16.04. The database server has 32GB RAM and 16
processors running Ubuntu 16.04.

#### 1

2	The user interface utilizes bootstrap version 3.4.1 for responsive sizing to different
3	format clients and browsers. Standard formatting, search, and filtering capability for
4	query tables is provided by the jQuery DataTables plugin. Highcharts is used for
5	generation of all plots. All queries and plots are produced dynamically from the
6	underlying database tables based on user query parameters. The URL is
7	https://gcregistry-explorer.stanford.edu/.

8

#### 9 RESULTS

#### 10 Participant population and demographics

11 From March 2011 to November 2021, 567 subjects enrolled in the study. For inclusion 12 in the study, all participants were required to identify their eligibility status on the 13 enrollment questionnaire. The majority reported a personal history of GC. Some 14 participants met multiple eligibility criteria (e.g., having GC and a family history of GC). Eligibility status with respect to sex, age, race, and ethnicity is depicted in Table 1. 15 16 Participants were predominantly female (63%), White (76%), and non-Hispanic (53%). 17 The median age of all participants was 51 years (range: 18-92 years). The median age 18 of participants with GC was 68 years. Some participants did not report their sex, race, 19 ethnicity, or other demographic details as these questions were optional on the 20 enrollment survey.

The participants' commonly reported medications included non-steroidal antiinflammatories, proton pump inhibitors and multivitamins. Common comorbidities
included high blood pressure, high cholesterol, and other cancers. Nearly a quarter of
participants with GC reported a history of gastric ulcers, gastroesophageal reflux
disease, or gastric polyps. *Helicobacter pylori* infection was reported by 16% of all
participants, and 20% of all participants with GC. Notably, a high proportion of
participants (65%) did not report or did not know their *H. pylori* status.

8

#### 9 Biospecimens

10 For this study, 164 participants donated biological specimens. We collected 111 saliva 11 samples and 89 tissue samples. At the time of publication, 41 tumor GCR samples had 12 undergone sequencing and with the results available on the GCR Genome Explorer 13 (Table 2). This subset of GCR archival samples underwent histopathological review, and most patients were diagnosed with gastric adenocarcinoma. Clinical and histologic 14 characteristics of the sequenced tumors are depicted in Table 2. Of the specimens 15 16 where Lauren's classification was reported (N=24), the majority were diffuse-type 17 cancers (N=16). The tumors were generally aggressive and poorly differentiated (59%). 18 The majority of tumors were from patients with metastatic disease (54%). With respect 19 to anatomic location, 10% of tumors arose from the cardia, 73% arose from the non-20 cardia stomach, and 17% arose from an unreported location.

21

#### 22 The GCR Genome Explorer and pilot data release

1 Public access is to the GCR Genome Explorer is available upon registration via the 2 following URL: https://gcregistry-explorer.stanford.edu/users/sign in. At time of 3 publication, the Genome Explorer contains data from the initial 41 sequenced tumors 4 through the GCR. In addition, genomic data from 443 TCGA gastric cancers and 185 5 TCGA esophageal cancers are available through the GCR Genome Explorer for cross-6 reference. Future releases will incorporate data from additional GCR tumor samples. A 7 representative image of the GCR Genome Explorer home page showing available data sets is depicted in Figure 1. 8

9

There are two tiers of results that are provided in the Genome Explorer. The first tier includes gene expression levels determined from RNA-seq, somatic copy number based on whole genome sequencing and somatic mutations derived from exome sequencing. The second tier leverages the first-tier results - different algorithms extrapolate characteristics of the clonal diversity (WGS, exome), cellular microenvironment (RNA-seq), microbiome content (RNA-seq), HLA genotypes (RNAseq) and putative neoantigens (WGS, exome, RNA-seq).

17

Cellular and microbiome features reflect the content of the local tumor
microenvironment. We extrapolated the cellular representation and microbial
populations using each tumor's RNA-seq data. The cell results were based on
processing with a deconvolution tool called CIBERSORTx.<sup>12</sup> The analysis
approximated the different cell types present in the tumor microenvironment.<sup>13</sup> For

those RNA-seq reads which did not align to the human genome, we used the Kraken2
program to determine if there were microbiome features that included bacterial
genera.<sup>14</sup>

4

5 Nonsynonymous mutations are a source of immunogenic peptides, called neoantigens, which are tumor-specific and not expressed in other normal cells. Tumor mutational 6 7 load, and more specifically, neoantigen load have been correlated with extent of T-cell 8 reactivity, response to checkpoint therapy, and prognosis.<sup>15-17</sup> Exome sequencing of these tumors allowed us to identify nonsynonymous mutations in the protein-coding 9 10 portions of genes in cancer patients and predict potential neoantigens. We identified 11 candidate neoantigens from the exome and RNA-seq data (Methods). A candidate 12 neoantigen fulfilled the following criteria: i) nonsynonymous somatic mutation, ii) 13 expressed in transcriptome data and iii) translated neopeptides with strong binding 14 affinity to patient's own major histocompatibility molecules. Using a combination of the 15 exome and transcriptome data and major histocompatibility complex genotypes, we 16 generated a list of potential neoantigens for each tumor.

17

#### 18 GCR Genome Explorer features

19 There are several ways of accessing results through the GCR Genome Explorer.

20 Options include general summaries of the results as well as specific queries. All image 21 files and tables are available for download. On the home page (**Figure 1**), the different 22 sample sets are displayed. For example, a user can query the GCR data set using the

1 "Explore Study" feature, which directs to a landing page with multiple tabs. In the 2 "Clinical Parameters" tab, the user will find visual and tabular representations of cohort 3 characteristics. The cohort can be gueried with regards to sex, race, ethnicity, cancer 4 site, age at diagnosis, body mass index, smoking history, cancer diagnosis, Lauren classification, and histologic differentiation. Under "Gene Summary," there are two 5 6 tables describing genes that are most frequently mutated or varied in copy number 7 within the GCR cohort. All cancer-associated genes are listed, along with genes mutated in >10% of samples, or genes with copy number variation in >10% of samples. 8 9 The percentage filter will change across cohorts due to rounding, and due to gene 10 ranking ties (i.e., multiple genes being mutated across the same total number of 11 samples). The tables include the gene name, cytoband, number of samples displaying 12 the mutation or copy number variation and ranked order of genes based on the 13 percentage of samples affected. In addition, we provide annotations for genes that are 14 cancer-associated, oncogenes or tumor-suppressor genes. Copy number summaries of 15 the pilot data set showed that many tumor samples had a high degree of genomic 16 instability. Whole genome sequencing revealed extensive changes in gene copy 17 number, with some amplified genes such as ERBB2 (i.e., HER2) being clinically 18 actionable. Amplifications of ERBB2 are an indication for the use of trastuzumab, a 19 therapeutic monoclonal antibody. Exome sequencing allowed us to detect specific 20 types of mutations. Each tumor sample contained unique combinations of missense 21 mutations, frame shift deletions, and in-frame deletions across various genes.

22

1 The tab "HLA Types" contains an overview and breakdown of specific HLA alleles for 20 2 patients in the GCR study. There are options to view the prevalence of specific HLA-A. 3 HLA-B, and HLA-C types within the study and within individual patients. In the next tab, 4 "Immune Cells," tumor-infiltrating immune cell types are identified and quantified in a bar 5 plot as a percentage of overall immune cells present within a patient's tumor. 6 Alternatively, users can view the immune cell presence represented as a heatmap. The 7 last tab, "Microbiome," contains an interactive bar chart displaying the major bacteria phylum found within each patient. Users can select and deselect phylum of interest for 8 9 a more granular, sample-specific view of the microbiome composition. 10 For more discrete summaries, users can employ the "Gene" query function to view data 11 12 for a specific gene. We cite an example using the *CDH1* tumor suppressor gene. 13 When searching all studies for CDH1 gene information, the user will find discrete 14 percentages and counts on its mutations, copy number variations, expression levels, and neoantigens within individual patients and across all cohorts (GCR and TCGA, 15 Figure 2). The CDH1 gene was mutated in 15.4% of GCR samples, 11% of TCGA-16 17 STAD samples and 1.1% of TCGA-ESCA samples. Nearly all samples in the GCR, 18 TCGA-STAD and TCGA-ESCA cohorts showed high expression of CDH1 (100%, 19 99.5%, and 99.4%).

20

The "Neoantigen" query gives the user the option to select an HLA allele type and viewa list of the candidate neoantigens. For example, the most prevalent HLA-A type in the

GCR-STCA study is A\*02:01 (18.9%). Selecting this HLA type in the Neoantigen query
produces a breakdown of all neoantigen candidates. They are described by their gene,
chromosomal location, amino acid change and binding strength/rank based on their
predicted properties in terms of MHC1 interaction (Figure 3). Like the "Gene" query,
the "Neoantigen" query function allows users to search data across multiple patient
cohorts.

7

8 Lastly, the "Patient" guery function provides users a way to see all available data for an 9 individual patient in the database (Figure 4). The patient's summary page includes a 10 count of mutations and copy number variations, their HLA-A, -B, and -C types, and their clinical characteristics. Like the "Explore Study" feature and the "Gene" query, the 11 12 "Patient" guery arranges datasets into tabs. The "Mutation" tab details the patient's 13 unique mutations with chromosome position, reference sequence and alteration, variant 14 type, and amino acid change. The "Copy Number" page lists genes with duplication and deletion events. In the "Gene Expression" tab, the user can view the expression 15 16 level of a gene through FPKM and qualitative values of high, medium, and low. The 17 gene list can be filtered to include only cancer-associated genes, oncogenes, or tumor 18 suppressor genes. The next tab contains putative neoantigens described by position, 19 amino acid change, and binding strength to a specific HLA allele. The final tab displays 20 the patient's microbiome, a taxonomy of microbial populations classified from kingdom 21 to genus.

22

#### 1 DISCUSSION

2 The GCR's goal is to catalyze research which will accelerate the prevention, detection, 3 and treatment of this malignancy. Patients with GC, a family history of GC in a first- or 4 second-degree relative, a germline *CDH1* mutation, or a combination of the three were 5 invited to enroll. Recruitment efforts were facilitated by making enrollment simple and 6 available to individuals who fulfill the eligibility criteria. Our national enrollment to 7 patients outside of the northern California region has enabled us to increase the 8 diversity of the participants. Through ongoing collaborations with other research groups 9 and healthcare centers across the globe, we are continuing to build out a robust 10 repository of clinical datasets, biospecimens and genomic data. The ongoing expansion 11 of the GCR will provide an even more diverse cohort and broaden the utility of this data 12 for research.

13

The GCR integrates epidemiology and genomic information from individuals at high risk 14 15 for GC or patients diagnosed with this malignancy. GCR Genome Explorer users can 16 conduct studies with this genomic information and take advantage of the sophisticated 17 bioinformatic tools used to process it. For example, a researcher can pose a question, "I found that *H. pylori* infected patients often have downstream functional mutation in the 18 19 gene *MUC1*; can this finding be replicated in another patient population?" The result is 20 provided efficiently among the GCR population as well as other studies such as the 21 TCGA. The information may be applicable to other cancer research studies as well.

22

1 While the pilot release of the Genome Explorer currently contains genomic data for 41 2 gastric tumors, our most recent enrollment will guadruple the number of samples 3 available for genomic studies. This influx of new participants is providing additional 4 clinical datasets such as pathology reports and other clinical metrics. On this expanded 5 cohort we are conducting additional genomic studies which will greatly increase the 6 overall number of tumors with genomic data. As the registry continues to accrue 7 participants and tumor samples, we anticipate that there will be improved statistical power for future studies. 8

9

10 Since the inception of the GCR, other groups have compiled clinical or genomic data on GC patients. The Stomach Cancer Pooling (StoP) Project is an international 11 12 consortium for epidemiological investigations into GC with 22 independent case-control studies participating.<sup>18</sup> Organized in 2015, the StoP consortium conducted multiple 13 studies to quantify the association between diet and lifestyle choices, 19-22 14 socioeconomic status,<sup>23</sup> and other factors<sup>24-26</sup> with the risk of GC. In 2015, researchers 15 identified discrete sources of epidemiological, clinicopathological, and molecular 16 biological information from GC researchers and treatment centers in China.<sup>27</sup> They 17 standardized the data to produce the Database of Human Gastric Cancer.<sup>27</sup> While 18 19 valuable, neither of these cited resources address the current need for an easily 20 accessible, comprehensive database that links detailed GC genomic analysis with 21 clinical data. What differentiates the GCR from these other efforts comes from its 22 integration of clinical and genomic data from GC patients, making it invaluable for future 23 translational studies. Moreover, there is no single genomic database of GCs based on

patients from the United States, a region with a lower prevalence of *H. pylori*-induced
cancer.<sup>28</sup>

3

For this iteration of the GCR's current registry we noted several limitations related to 4 5 patient reporting. First, participants in the GCR completed varying degrees of their questionnaires. All participants who completed our online consent form were included 6 7 in the final cohort. However, some participants did not complete every field of their 8 guestionnaire. Second, not all participants donated biological samples to the study. 9 Some patients elected not to contribute. This issue was exacerbated by the COVID-19 10 pandemic during which many patients were unable to donate saliva samples given the 11 biohazard issues. Third, using an online self-enrollment method limits participation to 12 individuals with access to the internet. One study suggested that older individuals and 13 those with lower self-rated internet ability are less likely to complete online enrollment processes.<sup>29</sup> Finally, our ongoing recruitment is expanding the overall diversity of the 14 cohort. 15

16

In summary, the GCR provides open access to a wealth of clinical and genomic information from a steadily expanding number of gastric tumors and individuals with hereditary predispositions for GC. We developed a framework for collating this information, analyzing tumor samples, and accessing this information via an online registry portal. Our results are openly accessible for independent discovery and validation studies.

1

## 2 ACKNOWLEDGEMENTS

- 3 The REDCap platform services at Stanford are subsidized by a) Stanford School of
- 4 Medicine Research Office, and b) the National Center for Research Resources and the
- 5 National Center for Advancing Translational Sciences, National Institutes of Health,
- 6 through grant UL1 TR001085.

#### 1 REFERENCES

- 2 1. Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A. Global cancer 3 statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 4 185 countries. CA Cancer J Clin. Nov 2018;68(6):394-424. doi:10.3322/caac.21492 5 2. Howlader N NA, Krapcho M, Miller D, Brest A, Yu M, Ruhl J, Tatalovich Z, Mariotto A, 6 Lewis DR, Chen HS, Feuer EJ, Cronin KA (eds). SEER Cancer Statistics Review, 1975-2017, 7 National Cancer Institute. Accessed May 1, 2020. https://seer.cancer.gov/csr/1975 2017/, 8 based on November 2019 SEER data submission, posted to the SEER web site, April 2020. 9 Surveillance Research Program NCI. SEER\*Explorer: An interactive website for SEER 3. 10 cancer statistics [Internet], Accessed May 1, 2020, [Cited 2020 Apr 15], Available from 11 https://seer.cancer.gov/explorer/. 12 4. Huang RJ, Choi AY, Truong CD, Yeh MM, Hwang JH. Diagnosis and Management of 13 Gastric Intestinal Metaplasia: Current Status and Future Directions. Gut Liver. Nov 15 14 2019;13(6):596-603. doi:10.5009/gnl19181 15 5. Zhou L, Catchpoole D. Spanning the genomics era: the vital role of a single institution 16 biorepository for childhood cancer research over a decade. Translational pediatrics. 17 2015;4(2):93-106. doi:10.3978/j.issn.2224-4336.2015.04.05 18 Liu A. Developing an institutional cancer biorepository for personalized medicine. Clinical 6. 19 Biochemistry. 2014/03/01/ 2014;47(4):293-299. 20 doi:https://doi.org/10.1016/j.clinbiochem.2013.12.015 21 McCarty CA, Chisholm RL, Chute CG, et al. The eMERGE Network: a consortium of 7. 22 biorepositories linked to electronic medical records data for conducting genomic studies. BMC 23 medical genomics. 2011;4(1):1-11. 24 Harris PA, Taylor R, Thielke R, Payne J, Gonzalez N, Conde JG. Research electronic 8. 25 data capture (REDCap)--a metadata-driven methodology and workflow process for providing 26 translational research informatics support. J Biomed Inform. Apr 2009;42(2):377-81. 27 doi:10.1016/j.jbj.2008.08.010 28 Harris PA, Taylor R, Minor BL, et al. The REDCap consortium: Building an international 9. 29 community of software platform partners. J Biomed Inform. Jul 2019;95:103208. 30 doi:10.1016/j.jbi.2019.103208 31 Xia LC, Van Hummelen P, Kubit M, et al. Whole genome analysis identifies the 10. 32 association of TP53 genomic deletions with lower survival in Stage III colorectal cancer. 33 Scientific Reports. 2020/03/19 2020;10(1):5009. doi:10.1038/s41598-020-61643-6 34 11. Institute NC. Genomic Data Commons Data Portal. October 29, 2021. 35 https://portal.qdc.cancer.gov/ 36 Newman AM, Steen CB, Liu CL, et al. Determining cell type abundance and expression 12. 37 from bulk tissues with digital cytometry. Nature Biotechnology. 2019/07/01 2019;37(7):773-782. 38 doi:10.1038/s41587-019-0114-2 39 Aran D, Hu Z, Butte AJ. xCell: digitally portraying the tissue cellular heterogeneity 13. 40 landscape. Genome Biol. Nov 15 2017;18(1):220. doi:10.1186/s13059-017-1349-1 41 Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. Genome 14. 42 Biol. 2019/11/28 2019;20(1):257. doi:10.1186/s13059-019-1891-0 43 Brown SD, Warren RL, Gibb EA, et al. Neo-antigens predicted by tumor genome meta-15. 44 analysis correlate with increased patient survival. Genome Research. 2014 2014;24(5):743-750. 45 doi:10.1101/gr.165985.113 Champiat S, Ferte C, Lebel-Binay S, Eggermont A, Soria JC. Exomics and 46 16. 47 immunogenics: Bridging mutational load and immune checkpoints efficacy. Oncoimmunology.
- 48 Jan 01 2014;3(1):e27817. doi:10.4161/onci.27817

- 1 17. Giannakis M, Mu XJ, Shukla SA, et al. Genomic Correlates of Immune-Cell Infiltrates in 2 Colorectal Carcinoma. *Cell Reports*. 2016 2016;15(4):857-865.
- 3 doi:10.1016/j.celrep.2016.03.075

4 18. Pelucchi C, Lunet N, Boccia S, et al. The stomach cancer pooling (StoP) project: study 5 design and presentation. *Eur J Cancer Prev*. Jan 2015;24(1):16-23.

- 6 doi:10.1097/CEJ.000000000000017
- 7 19. Rota M, Pelucchi C, Bertuccio P, et al. Alcohol consumption and gastric cancer risk-A
  8 pooled analysis within the StoP project consortium. *Int J Cancer*. Nov 15 2017;141(10):1950-
- 9 1962. doi:10.1002/ijc.30891
- 10 20. Praud D, Rota M, Pelucchi C, et al. Cigarette smoking and gastric cancer in the Stomach
- 11 Cancer Pooling (StoP) Project. *Eur J Cancer Prev*. Mar 2018;27(2):124-133.
- 12 doi:10.1097/CEJ.000000000000290
- 13 21. Bertuccio P, Alicandro G, Rota M, et al. Citrus fruit intake and gastric cancer: The
- stomach cancer pooling (StoP) project consortium. *Int J Cancer*. Jun 15 2019;144(12):29362944. doi:10.1002/ijc.32046
- 16 22. Ferro A, Rosato V, Rota M, et al. Meat intake and risk of gastric cancer in the Stomach 17 cancer Pooling (StoP) project. *Int J Cancer*. Oct 4 2019;doi:10.1002/ijc.32707
- 18 23. Rota M, Alicandro G, Pelucchi C, et al. Education and gastric cancer risk-An individual
- 19 participant data meta-analysis in the StoP project consortium. Int J Cancer. Feb 1
- 20 2020;146(3):671-681. doi:10.1002/ijc.32298
- 21 24. Ferro A, Morais S, Pelucchi C, et al. Smoking and Helicobacter pylori infection: an
- individual participant pooled analysis (Stomach Cancer Pooling- StoP Project). *Eur J Cancer Prev.* Sep 2019;28(5):390-396. doi:10.1097/CEJ.000000000000471
- 24 25. Ferro A, Morais S, Pelucchi C, et al. Sex differences in the prevalence of Helicobacter
  25 pylori infection: an individual participant data pooled analysis (StoP Project). *Eur J Gastroenterol*26 *Hepatol.* May 2019;31(5):593-598. doi:10.1097/MEG.00000000001389
- 27 26. Shah SC, Boffetta P, Johnson KC, et al. Occupational exposures and odds of gastric
- cancer: a StoP project consortium pooled analysis. *Int J Epidemiol*. Jan 21
- 29 2020;doi:10.1093/ije/dyz263
- 30 27. Wang C, Zhang J, Cai M, et al. DBGC: A Database of Human Gastric Cancer. *PLoS* 31 *One*. 2015;10(11):e0142591. doi:10.1371/journal.pone.0142591
- 28. Nguyen TH, Mallepally N, Hammad T, et al. Prevalence of Helicobacter pylori Positive Non-cardia Gastric Adenocarcinoma Is Low and Decreasing in a US Population. *Dig Dis Sci.*
- 34 Nov 14 2019;doi:10.1007/s10620-019-05955-2
- 35 29. Buis LR, Janney AW, Hess ML, Culver SA, Richardson CR. Barriers encountered during
- 36 enrollment in an internet-mediated randomized controlled trial. *Trials*. Aug 23 2009;10:76.
- 37 doi:10.1186/1745-6215-10-76
- 38

# 1 DATA ACCESSIBILITY 2

- 3 Primary data is available on the Gastric Genome Explorer website (https://gcregistry-
- 4 explorer.stanford.edu/). Other data generated in this study are available within the
- 5 article and its supplementary data files.

## Table 1. Cohort characteristics by eligibility status from 2011-2021. As participants could fall into more than one eligibility category, some participants appear in multiple

columns.

	All Participants (N = 567)	Gastric Cancer (N = 366)	Family History (N = 206)	CDH1 Mutation (N = 78)				
	Frequency (%)							
Sex								
Female	357 (63%)	195 (53%)	160 (78%)	66 (85%)				
Male	208 (37%)	170 (46%)	45 (22%)	12 (15%)				
Unknown	2 (<1%)	1 (<1%)	1 (<1%)	0 (0%)				
Age (years)								
<40	144 (25%)	61 (17%)	77 (37%)	28 (36%)				
40-49	114 (20%)	76 (21%)	41 (20%)	17 (22%)				
50-59	125 (22%)	92 (25%)	40 (19%)	17 (22%)				
60-69	107 (19%)	77 (21%)	28 (14%)	13 (17%)				
70-79	6(1%)	6 (2%)	0 (0%)	0 (0%)				
≥80	11 (2%)	8 (2%)	4 (2%)	0 (0%)				
Unknown	60 (11%)	46 (13%)	16 (8%)	3 (4%)				
Race/Ethnicity								
White	432 (76%)	271 (74%)	158 (77%)	69 (88%)				
Black	24 (4%)	20 (5%)	6 (3%)	2 (3%)				
Native American	1 (<1%)	1 (<1%)	0 (0%)	0 (0%)				
Asian	43 (8%)	34 (9%)	12 (6%)	4 (5%)				
Pacific Islander	2 (<1%)	1 (<1%)	1 (<1%)	0 (0%)				
Other	45 (8%)	25 (7%)	23 (11%)	3 (4%)				
Missing	20 (4%)	14 (4%)	6 (3%)	0 (0%)				
Ethnicity								
Hispanic	83 (15%)	52 (14%)	35 (17%)	6 (8%)				
Non-Hispanic	443 (78%)	287 (78%)	154 (75%)	67 (86%)				
Unknown	41 (7%)	27 (7%)	17 (8%)	5 (6%)				
Comorbidities								
Diabetes	69 (12%)	57 (16%)	12 (6%)	8 (10%)				
Hyperlipidemia	129 (23%)	85 (23%)	48 (23%)	18 (23%)				
Hypertension	134 (24%)	99 (27%)	42 (20%)	19 (24%)				
Helicobacter pylori								
Tested positive	91 (16%)	73 (20%)	22 (11%)	6 (8%)				
Tested negative	108 (19%)	70 (19%)	37 (18%)	17 (22%)				
Did not test	368 (65%)	223 (61%)	147 (71%)	55 (71%)				
Risk Factors								
Epstein-Barr virus	11 (2%)	5(1%)	5 (2%)	2 (3%)				
Gastric ulcer	112 (20%)	83 (23%)	34 (17%)	14 (18%)				
Gastric polyps	40 (7%)	22 (6%)	22 (11%)	3 (4%)				
Gastroesophageal reflux disease	131 (23%)	77 (21%)	56 (27%)	11 (14%)				

## 1 Table 2. Overview of gastric tumor tissues in the GCR Genome Explorer (N=41).

	Frequency (%)
Sex	
Male	24 (59%)
Female	16 (39%)
Not available	1 (2%)
Histology	(_,,,)
Gastric adenocarcinoma	37 (90%)
GIST	1 (2%)
Not available	2 (5%)
Histological Subtype	_ (***)
Intestinal	4 (10%)
Diffuse	16 (39%)
Mixed	4 (10%)
Not available	17 (41%)
Histological differentiation	(,)
Well	1 (2%)
Moderate	3 (7%)
Moderately to poorly	3 (7%)
Poorly	24 (59%)
Not available	10 (24%)
AJCC tumor pathology	
T1	6 (15%)
T2	0 (0%)
ТЗ	12 (29%)
T4	7 (17%)
Not available	16 (39%)
N0	9 (22%)
N1	3 (7%)
N2	4 (10%)
N3	7 (17%)
Not available	18 (44%)
MO	1 (2%)
MX	22 (54%)
Not available	18 (44%)
Lymph node involvement	
Positive	13 (32%)
Negative	12 (29%)
Not available	16 (39%)
Adjuvant radiation	
Yes	14 (34%)
No	20 (49%)
Not available	7 (17%)
Tumor anatomic site	
Gastroesophageal junction	5 (12%)
Cardia	4 (10%)
Pylorus	2 (5%)
Fundus	3 (7%)
Body	8 (20%)
Antrum	5 (12%)
Entire stomach	7 (17%)
Not available	7 (17%)

## 1 FIGURES

- 2 Figure 1. The GCR Genome Explorer home page. The portal offers queries of gene,
- 3 neoantigen and patient datasets from one or multiple cohorts: the GC Registry (GCR),
- 4 the Cancer Genome Atlas (TCGA) esophageal carcinoma (ESCA) study, and the TCGA
- 5 stomach adenocarcinoma (STAD) study. The "Explore Study" feature enables the user
- 6 to see the clinical and genomic characteristics of a single cohort.

elect study/studies to query:				Example queries
Select/Deselect All				<ul> <li>Explore Study:</li> <li>What is the age distribution of patients in the GCR STAD study</li> <li>How many GCR STAD patients have an <i>FGFR2</i> amplification?</li> </ul>
Stomach cancer		41 samples	Explore Study	<ul> <li>What immune cell types are enriched in GCR STAD patients?</li> <li>What is the proportion of bacterial types in the microbiomes of GCR STAD patients?</li> </ul>
TCGA				Gene Query:
Esophageal carcinoma		185 samples	Explore Study	<ul> <li>How many TCGA STAD patients have a missense mutation in TP53?</li> </ul>
Stomach adenocarcinoma		443 samples	Explore Study	How does the proportion of <i>KRAS</i> mutations compare between GCR and TCGA STAD studies?
uery the study/studies by:				Neoantigen/HLA Query:     Are there any GCR STAD patients with neoantigens that have a
				strong binding affinity to their own HLA type?

- 1 **Figure 2.** The GCR Genome-Explorer output of a "Gene" query across all study
- 2 cohorts for CDH1. The "Summary" page displays counts and percentages of samples
- 3 with CDH1 mutations, copy number variations, and expression levels. The "Mutation,"
- 4 "Copy Number," "Expression," and "Neoantigen" tabs contain more detailed gene
- 5 alterations among the studies' patients.



- 1 Figure 3. Example of the "Neoantigen" query for HLA type A\*02:01 within the GCR-
- 2 STCA cohort. The table describes specific gene mutations that lead to the production of
- 3 a neoantigen and the predicted binding strength of the neoantigen to HLA allele
- 4 A\*02:01.

GCR ge	enome explore	ər							Pro	ect info -	GQ
Necantigen/HLA Query HLA Type: A*02.01 Cancer(Study): STCAIGCR Number of samples in queried studies: 41											
how 50 $\sim$	entries									Se	arch:
Study 斗	Cancer 1	Patient 1	Gene ⊥↑	ChrPos 11	Ref/Alt 1	AAChange 11	Mut TPM 1	HLA Allele	Neoantigen 1	Binding Strength	Binding Rank
GCR	STCA	P05947	PLEKHH1	chr14_67578543	C/T	p.Leu921Phe	0.331064	A*02:01	SLMQCWQFL	weak	1.491
GCR	STCA	P05936	CYP2C19	chr10_94781864	C/T	p.Thr229lle	0.989242	A*02:01	TIIDYFPGI	strong	0.0777
GCR	STCA	P05936	MS4A7	chr11_60389457	C/T	p.Ala136Val	0.168328	A*02:01	GLFLLVDSM	weak	0.8282
GCR	STCA	P05936	MS4A7	chr11_60389457	C/T	p.Ala136Val	0.168328	A*02:01	FLLVDSMVA	strong	0.4217
GCR	STCA	P05936	MS4A7	chr11_60389457	С/Т	p.Ala136Val	0.168328	A*02:01	LLVDSMVAL	strong	0.0578
GCR	STCA	P05936	ESAM	chr11_124756679	G/A	p.Pro105Ser	43.202485	A*02:01	SMSSRNLSL	weak	0.9322
GCR	STCA	P05936	PAN3	chr13_28261456	С/Т	p.Pro470Leu	0.431154	A*02:01	AQIDQADML	weak	1.2157
GCR	STCA	P05936	PAN3	chr13_28261456	С/Т	p.Pro470Leu	0.431154	A*02:01	DMLAVPTEV	weak	1.1397
GCR	STCA	P05936	PPP1R3E	chr14_23302525	T/C	p.lle18Val	4.459379	A*02:01	FVAALTERA	weak	1.5433
GCR	STCA	P05936	HERC2	chr15_28222082	C/T	p.Met1866lle	24.856053	A*02:01	AIMKIGTRV	strong	0.4498
GCR	STCA	P05936	CAMKK1	chr17_3872605	A/G	p.lle358Thr	1.41339	A*02:01	FIDDFTLAL	strong	0.0143
GCR	STCA	P05936	TRPM4	chr19_49190251	C/T	p.Thr688lle	0.6267	A*02:01	STIPIWALV	weak	1.031
GCR	STCA	P05936	GEN1	chr2_17781815	C/T	p.Pro868Leu	0.103536	A*02:01	FLDSTKSSL	strong	0.1273
GCR	STCA	P05936	RGPD8	chr2_112380827	A/C	p.lle1686Met	0.933747	A*02:01	VLMEQMKLL	strong	0.0902
GCR	STCA	P05936	PPP1R7	chr2_241166343	C/T	p.Arg241Trp	0.459748	A*02:01	MQSNWLTKI	weak	1.7772
GCR	STCA	P05936	CRBN	chr3_3151040	G/A	p.Ala385Val	0.295263	A*02:01	WFPGYVWTV	weak	1.8836
GCR	STCA	P05936	RPSA	chr3_39411972	С/Т	p.Ala235Val	9.181129	A*02:01	FQGEWTAPV	strong	0.2904
GCR	STCA	P05936	DNAJC13	chr3_132538235	C/T	p.Pro2229Ser	0.810564	A*02:01	VMSNLSPPV	strong	0.4443
GCR	STCA	P05936	FBXW7	chr4_152328358	C/A	p.Gly423Val	3.342483	A*02:01	TLVGHTGVV	weak	1.5289
GCR	STCA	P05936	CMTR1	chr6_37462872	C/T	p.Arg457Trp	1.359764	A*02:01	GIDDVWDYL	strong	0.4193
GCR	STCA	P05936	CMTR1	chr6_37462872	С/Т	p.Arg457Trp	1.359764	A*02:01	DVWDYLFAV	strong	0.4403
GCR	STCA	P05936	FAM120B	chr6_170318205	C/T	p.Ser272Leu	0.729989	A*02:01	ILAVSDHIL	weak	0.6055
GCR	STCA	P05936	FAM120B	chr6_170318205	C/T	p.Ser272Leu	0.729989	A*02:01	AVSDHILKV	strong	0.1121
GCR	STCA	P05936	ECPAS	chr9_111389716	G/A	p.Ala1096Val	0.202546	A*02:01	AVEGENVIA	weak	0.7474
GCB	STCA	P06872	CYP46A1	chr14 99706775	T/C	p.lle191Thr	1.577008	A*02:01	TLAKAAEGM	weak	0.926

- 1 Figure 4. "Patient" query for P04906 from the GCR study. Multiple tabs present the
- 2 entirety of clinical, genomic, and cellular data for this individual.

