1  **TITLE**: Clinical interpretation of integrative molecular profiles to guide precision cancer medicine

2

3  **AUTHORS**: Brendan Reardon[1,2], Nathaniel D Moore[1,2,3,4,5], Nicholas Moore[1,2,6], Eric
4  Kofman[1,2,7,8], Saud Aldubayan[1,2,9,10], Alexander Cheung[1,2,11], Jake Conway[1,2,12], Haitham
5  Elmarakeby[1,2,13], Alma Imamovic[2,14], Sophia C. Kamran[2,15], Tanya Keenan[1,2], Daniel Keliher[1,2,16],
6  David J Konieczkowski[2,17,18,19], David Liu[1,2], Kent Mouw[2,6,17], Jihye Park[1,2], Natalie Vokes[1,2,20],
7  Felix Dietlein[1,2], Eliezer M Van Allen[1,2]*

8

9  **AFFILIATIONS**: Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA,
10  USA[1]; Broad Institute of MIT and Harvard, Cambridge, MA, USA[2]; Indiana University School of
11  Medicine, Indianapolis, IN, USA[3]; Howard Hughes Medical Institute, Chevy Chase, MD, USA[4];
12  Department of Internal Medicine, University of Cincinnati, Cincinnati, Ohio, USA[5]; Harvard
13  Medical School, Harvard University, Boston, MA, USA[6]; Department of Cellular and Molecular
14  Medicine, University of California, San Diego, La Jolla, CA, USA[7]; Institute for Genomic
15  Medicine, University of California, San Diego, La Jolla, CA, USA[8]; Division of Genetics, Brigham
16  and Women's Hospital, Boston, MA, USA[9]; Division of Genetics, Brigham and Women's
17  Hospital, Boston, MA, USA[9]; College of Medicine, King Saud bin Abdulaziz University for Health
18  Sciences, Riyadh, Saudi Arabia[10]; Grossman School of Medicine, New York University, New
19  York, NY, USA[11]; Division of Medical Sciences, Harvard University, Boston, MA, USA[12];
20  Department of System and Computer Engineering, Al-Azhar University, Cairo, Egypt[13];
21  Department of Pediatric Oncology, Dana-Farber Cancer Institute, Harvard Medical School,
22  Boston, MA, USA[14]; Department of Radiation Oncology, Massachusetts General Hospital,
23  Harvard Medical School, Boston, MA, USA[15]; Department of Mathematics, Tufts University,
24  Medford, MA, USA[16]; Department of Radiation Oncology, Dana-Farber Cancer Institute &
25  Brigham and Women's Hospital, Boston, MA[17]; Harvard Radiation Oncology Program,
26  Massachusetts General Hospital, Boston, MA, USA[18]; Department of Radiation Oncology, The
27  Ohio State University Comprehensive Cancer Center - Arthur G. James Cancer Hospital and
28  Richard J Solove Research Institute, Columbus, OH, USA[19]; Department of Thoracic / Head
29  and Neck Oncology, MD Anderson Cancer Center, Houston, TX, USA[20]

30

31  **CORRESPONDING AUTHOR**: Eliezer M Van Allen (EliezerM_VanAllen@dfci.harvard.edu)

32

33  **WORD COUNT**: 3650

## ABSTRACT

Individual tumor molecular profiling is routinely used to detect single gene-variant ("first-order") genomic alterations that may inform therapeutic actions -- for instance, a tumor with a *BRAF* p.V600E variant might be considered for RAF/MEK inhibitor therapy. Interactions between such first-order events (e.g., somatic-germline) and global molecular features (e.g. mutational signatures) are increasingly associated with clinical outcomes, but these "second order" alterations are not yet generally accounted for in clinical interpretation algorithms and knowledge bases. Here, we introduce the Molecular Oncology Almanac (MOAlmanac), a clinical interpretation algorithm paired with a novel underlying knowledge base to enable integrative interpretation of genomic and transcriptional cancer data for point-of-care treatment decision-making and translational hypothesis generation. We compared MOAlmanac to first-order interpretation methodology in multiple retrospective patient cohorts and observed that the inclusion of preclinical and inferential evidence as well as second-order molecular features increased the number of nominated clinical hypotheses. MOAlmanac also performed matchmaking between patient molecular profiles and cancer cell lines to further expand individualized clinical actionability. When applied to a prospective precision oncology trial cohort, MOAlmanac nominated a median of two therapies per patient and identified therapeutic strategies administered in 46% of patient profiles. Overall, we present a novel computational method to perform integrative clinical interpretation of individualized molecular profiles. MOAlmanc increases clinical actionability over conventional approaches by considering second-order molecular features and additional evidence sources, and is available as an open-source framework.

# INTRODUCTION

Targeted panels or whole-exome sequencing now routinely inform the clinical care of oncology patients[1]. The resulting collections of patient-specific cancer genome alterations are valuable resources in the advancement of precision medicine. However, the growing quantity and complexity of potentially actionable genomic alterations available for each patient limit the ability of any individual clinician or researcher to interpret them. This challenge necessitated the creation of clinical interpretation algorithms to computationally prioritize large sets of patient-specific alterations by clinical and biological relevance, as well as exposed the need to pair these interpretation algorithms with up-to-date knowledge bases that link molecular alterations to relevant clinical actions.

Clinical decision-making in precision oncology commonly emphasize "first-order" relationships -- pairing individual somatic variants, copy number alterations, pathogenic germline variants, or fusions with specific clinical actions such as use of *BRAF* p.V600E and RAF/MEK inhibition -- based on FDA approvals and other clinical evidence[2-7]. While these efforts have been highly fruitful, they also have certain limitations. Many academic and commercially available targeted panels focus primarily on somatic variants and copy number alterations; often, they do not sequence associated germline tissue or comprehensively assess fusions[1]. Yet pathogenic germline variants impact cancer risk and can also modify clinical interpretation of secondary somatic events in the same gene or of genome-wide mutational signatures, e.g. DNA repair[8,9]. Similarity, the approval of TRK inhibitors for patients with any solid tumor harboring *NTRK* fusions and other biological insights gained from somatic variants that can be identified from RNA may warrant expanding routine clinical sequencing to jointly evaluate a patient's genomic and transcriptional data[10,11]. In addition, the ongoing characterization of the cancer genome has revealed the importance of considering these first-order events in tandem as well as "second-order" molecular features -- genomic processes such as microsatellite instability and tumor mutational burden that are global rather than limited to individual gene(s). Such processes have also been associated with clinical phenotypes, such as COSMIC Signature 6 correlating with mismatch repair deficiency (MMR) and microsatellite instability (MSI) linked to cancer immunotherapy response[12]. Lastly, even with the consideration of these additional features and second-order relationships, some patients may be variant-negative and thus may not qualify for genomically guided treatment. To address this challenge, multiple efforts have demonstrated that preclinical cell line models can also inform treatment selection, but such approaches are constrained by both the limited molecular diversity of cancer cell lines and computational difficulty in matchmaking, to identify which models are most representative of an individual patient's tumor[13-17].

To maximize interpretability of integrative molecular profiling for point-of-care treatment decision-making and translational hypothesis generation, new methodologies are needed to leverage both first- and second-order molecular alterations, relationships between multiple co-occurring events, and the full spectrum of both clinical and preclinical evidence. Here, we introduce Molecular Oncology Almanac (MOAlmanac), a clinical interpretation algorithm paired with an alteration-action database (Figure 1) that operates on germline, somatic, and

99   transcriptional data in tandem from individual patients. MOAlmanac expands the scope of
100  considered molecular alterations beyond somatic variants and copy number alterations to
101  include fusions, germline variants, and concordance between events across feature types. In
102  addition, MOAlmanac considers global "second-order" molecular features and introduces a
103  patient-to-cell line matchmaking module to leverage cell line profiling to nominate additional
104  genomic features potentially associated with therapeutic sensitivity. MOAlmanac is provided in a
105  cloud-based framework and delivers reports at the level of the individual patient. By integrating
106  diverse data sources with higher-order interpretation, MOAlmanac expands the landscape of
107  clinical actionability to facilitate point-of-care decision making and to advance precision cancer
108  medicine.


# RESULTS

109

## Developing an integrated interpretation framework

110

111  Molecular Oncology Almanac is a clinical interpretation method that evaluates individual patient
112  molecular profiles to facilitate precision oncology (Figure 1a). Individual genomic events are
113  annotated and sorted to identify those that are both highly associated with cancer and
114  associated with treatment response or prognosis. First, features are prioritized based on an
115  association between the involved genes and cancer in several data sources; in order:
116  MOAlmanac's database (described below), Cancer Hotspots, 3D Cancer Hotspots, Cancer
117  Gene Census (CGC), Molecular Signatures Database, and Catalogue of Somatic Mutations in
118  Cancer (COSMIC) (Methods, Supplementary Figure 1a)[18–23]. Next, molecular features are
119  further prioritized based on associations between specific alterations and each data source. For
120  instance, *KRAS* p.G12A ranks higher than *KRAS* p.I36M as both protein changes are reported
121  as 3D hotspots but only p.G12A matches to Cancer Hotspots.
122
123  The clinical relevance of each cancer-associated molecular feature is further assessed based
124  on an underlying custom knowledge base, which contains 722 assertions relating molecular
125  features to therapeutic sensitivity, resistance, and prognosis based on published literature and
126  guidelines. This resource evolved from our prior actionability database (Tumor Alterations
127  Relevant for GEnomics-driven Therapy (TARGET)), which represented entries as genes and
128  data types[2] (Figure 1b, Methods). In contrast, MOAlmanac defines molecular features broadly to
129  encompass the varying types of alterations backed by cited evidence. For example,
130  MOAlmanac is capable of recording information regarding specific singleton features (e.g.
131  *BRAF* p.V600E) but also more general event classes (such as the presence of an *ALK* fusion
132  without regard to the fusion partner). Relationships between molecular features and treatment
133  response are annotated for targeted therapies (415 assertions), immunotherapies (48),
134  chemotherapies (40), radiation therapy (15), hormonal treatments (7), and combination
135  therapies (11) (Figure 1c, Methods). Individual genomic events that match catalogued features
136  are labeled by the specificity of the underlying event and match completeness. For example,
137  exact matches to fully defined features, such as *BCR-ABL1*, are labeled as "Putatively
138  Actionable"; partial matches within a feature type are labeled as "Investigate Actionability", such
139  as an *ATM* missense variant matching to a catalogued *ATM* nonsense variant; and events
140  whose gene appears in the database under a different data type are highlighted as "Biologically

141  Relevant" but not associated with a clinical assertion, e.g. a *CDKN2A* somatic variant matching
142  to *CDKN2A* copy number deletions. These assertions are derived from numerous evidence
143  sources in accordance with existing frameworks[3–5,24], including: FDA approvals (FDA-approved),
144  clinical guidelines (Guideline), results from prospective clinical trials (Clinical trial), results from
145  human studies other than a clinical trial (Clinical evidence), findings from cancer cell lines or
146  animal models (Preclinical), or inferences from mathematical models or associations between
147  molecular features (Inferential) (Figure 1c, Methods).
148
149  MOAlmanac also characterizes individual features in concert with each other and second-order
150  genomic events. For each MOAlmanac gene, events across all feature types are reported
151  together to elucidate contributions from distinct types of genomic events. Somatic variants in a
152  given gene will increase in priority if either a truncating or a pathogenic or likely pathogenic
153  (according to ClinVar) germline variant appears in the same gene or if the somatic variant is
154  observed with sufficient power in validation sequencing, if provided[24,25]. Both COSMIC
155  mutational signature contributions and tumor mutational burden (TMB) are calculated and
156  variants related to microsatellite instability are highlighted. Tumor ontology is mapped with
157  Oncotree. Tumor purity, ploidy, whole-genome doubling, and microsatellite stability status are
158  also accepted for reporting and evaluation. All nominated clinical associations are reported in a
159  web-based actionability report (Methods).


160  **Evaluating expanded molecular profiling and actionability in two retrospective cohorts**

161  We first evaluated MOAlmanac relative to our prior established whole-exome sequencing
162  (WES) first-order interpretation framework (PHIAL with TARGET), which considers somatic
163  variants and copy number alterations[2]. WES and RNA-sequencing (RNA-seq) data were
164  acquired for 110 previously published metastatic melanomas (n = 44 with RNA)[26] and 150
165  patients with metastatic castration-resistant prostate cancers (mCRPC, n = 149 with RNA)[27]. All
166  samples were analyzed to call somatic variants, germline variants, and copy number alterations
167  from WES and somatic variants and fusions from RNA-seq (Methods).
168
169  We compared how often the two methods observed a clinically relevant event associated with
170  therapeutic sensitivity, resistance, or prognosis when only somatic variants and copy number
171  alterations were considered. Furthermore, we characterized only well-established relationships
172  by restricting our analysis to assertions curated from FDA approvals, clinical guidelines, clinical
173  trials, or clinical evidence. MOAlmanac identified 312 such putatively actionable events from
174  191 patients (73 melanoma, 118 mCRPC), 218 (69.87%) of which were flagged by PHIAL for
175  clinical relevance. For example, the most commonly flagged features were *BRAF* p.V600E (39
176  patients), *MET* amplification (9), and *PTEN* deletion (9) for metastatic melanomas and *AR*
177  amplifications (82), *PTEN* deletions (40), and *RB1* deletions (21) in mCRPC. When "Investigate
178  Actionability" variants were included, an additional 54 patients (20.8% of cohort) harbored a
179  potentially clinically relevant variant, such as *NRAS* p.Q61K (10, melanoma) with associated
180  sensitivity to selumetinib, 31 of which were also highlighted by PHIAL. PHIAL identified 0 events
181  as Putatively Actionable and 113 as Investigate Actionability which were not highlighted by
182  MOAlmanac; however, all genes associated with these events were not migrated to

183    MOAlmanac from TARGET for reasons such as insufficient evidence of clinical relevance
184    (Methods).
185
186    Next, while still limiting our analysis to somatic variants and copy number alterations, we
187    investigated how the inclusion of preclinical and inferential  evidence sources affected
188    identification of potentially actionable results. On the basis of preclinical evidence, 120 such
189    genomic events from 107 patients were identified -- for example, *PTEN* deletions and sensitivity
190    to everolimus or AZD8186, 86 (71.7%) of which were also highlighted by PHIAL. Inferential
191    evidence highlighted 19 additional putatively actionable copy number alterations from 19
192    patients, most prominently *CCND1* amplifications for reported sensitivity to palbociclib (n=15).
193    Thus, using all catalogued evidence, MOAlmanac noted 1175 somatic variants and copy
194    number alterations as Putatively Actionable or Investigate Actionability across 249 patients (109
195    melanoma, 140 CRPC). Of these events, PHIAL highlighted 73 (6.2%) as Putatively Actionable,
196    352 (30%) as Investigate Actionability, and 369 (31.4%) as Biologically Relevant.
197
198    We then evaluated whether an expanded set of molecular features (including germline variants
199    and fusions as additional first-order features and tumor mutational burden, mutational
200    signatures, and aneuploidy as second-order features, none of which are handled by PHIAL,
201    could further broaden the actionability landscape for individual patients (Figure 2b). Pathogenic
202    and likely pathogenic germline variants highlighted 10 additional clinically relevant molecular
203    features across 10 different samples (0 melanoma, 10 mCRPC), six of which were *BRCA1/2*
204    variants. MOAlmanc identified 127 clinically relevant fusions across 82 patients; ten mCRPC
205    tumors harbored no putatively actionable somatic variants or copy number alterations but did
206    contain *TMPRS22-ERG*. Regarding second-order molecular features, elevated TMB was noted
207    for 43 patients with metastatic melanoma and 4 with mCRPC (Methods), clinically relevant
208    mutational signatures were observed in 40 molecular profiles, and whole-genome doubling,
209    which has been associated with poor prognosis, was observed in 137 profiles[28]. In some of
210    these cases, combinations of these features were particularly relevant when present in tandem.
211    For example, a pathogenic *BRCA2* variant, p.S1882*, was observed in one patient along with a
212    39% mutational signature attribution to COSMIC Signature 3, both of which may suggest
213    homologous recombination repair deficiency (HRD) and sensitivity to PARP inhibition[29–31]. By
214    considering these feature types, MOAlmanac identified an additional 397 clinically relevant
215    molecular features in 214 patients, resulting in 258 patients with at least one event associated
216    with therapeutic sensitivity, resistance, or prognosis.
217
218    Focusing specifically on therapeutic sensitivity, such consideration of an extended set of feature
219    types and additional evidence sources provided otherwise variant-negative patients with clinical
220    hypotheses (Figure 2c, Supplementary Table 1). FDA approved or clinical guideline
221    associations resulted in a highlighted therapy for 175 of 260 patients (75 and 100 for melanoma
222    and CRPC, respectively); 11 patients obtained a therapeutic hypothesis from feature types other
223    than somatic variants and copy number alterations, such as elevated TMB (2 patients) or *NTRK*
224    fusions (1). Inclusion of preclinical and inferential evidence sources further decreased the
225    number of variant-negative patients from 85 to 11 (41 preclinical, inferential); for example
226    *CDKN2A/B* deletions and sensitivity to EPZ015666 (6).

227

228  In total, MOAlmanac found at least one clinically relevant feature in 100% and 98.6% of
229  metastatic melanoma and mCRPC profiles, using evidence ranging from FDA approvals to
230  inferential relationships and both first- and second-order molecular features (Figure 2a, 2b). In
231  comparison, PHIAL identified such somatic variants and copy number alterations in 92.7% and
232  89.3% of metastatic melanoma and mCRPC profiles, respectively. Thus, the inclusion of
233  additional feature types and evidence for clinical interpretation provided patients with an
234  expanded set of clinical hypotheses.


235  **Leveraging preclinical models for clinical actionability**

236  We next investigated whether preclinical data from high-throughput therapeutic screens of
237  cancer cell lines could further inform clinical interpretation within the MOAlmanc methodology.
238  We identified 452 solid tumor cell lines from the Cancer Cell Line Encyclopedia (CCLE) and
239  Sanger Institute's Genomics of Drug Sensitivity in Cancer (GDSC) that had available data on
240  nucleotide variants, copy number alterations, fusions, and drug sensitivity (Methods)[32,33]. Of
241  MOAlmanac's 124 catalogued therapies, 44 were represented in the current GDSC2 dataset
242  and 15 additional therapies were represented only in the older GDSC1 dataset. These 44
243  therapies are involved in 159 catalogued assertions between genomic alterations and
244  therapeutic sensitivity, for each MOAlmanac evaluates sensitivity for wild-type cell lines vs those
245  harboring the corresponding or related alterations. For example, in the case of the catalogued
246  preclinical relationship between *PIK3CA* p.H1047R and sensitivity to pictilisib, MOAlmanac
247  reports sensitivity for wild-type cell lines versus those harboring any genomic alteration in
248  *PIK3CA*, any nonsynonymous variant in *PIK3CA*, any missense variant in the gene, and those
249  specifically with the p.H1047R variant (Supplementary Figure 3a). Across all evaluable
250  relationships asserting sensitivity, 12 therapies showed a significant difference in IC50 between
251  wild type and mutant cell lines (Supplementary Table 2, Methods). Thus, high-throughput
252  therapeutic screens of cancer cell lines are used as an orthogonal axis of evidence to evaluate
253  clinically relevant relationships nominated by MOAlmanac.
254

255  The above approach simplistically compares sensitivity between cell lines that do or do not
256  share a single specific molecular feature. A potential limitation of this approach is that it includes
257  cell lines that share the index feature but are otherwise genomically highly dissimilar and
258  therefore whose overall biological relevance to the underlying patient sample may be
259  questionable. Therefore, we were motivated to identify cancer cell lines that shared more
260  extensive similarities in their molecular profiles and investigate whether such "patient-to-cell line
261  matchmaking" could identify additional potential therapeutic sensitivities. Previous approaches
262  have evaluated genomic similarity based on shared mutated genes that are weighted by their
263  recurrence in TCGA[15,16]; however, we chose to assess models based on shared therapeutic
264  sensitivity independent of histology-specific priors. We evaluated several models on cell lines
265  using a hold-one-out approach (Methods). For each cell line, we determined whether its nearest
266  neighbor shared drug sensitivity to any GDSC therapy (Figure 3a, Methods). Similarity Network
267  Fusion applied to nucleotide variants, copy number alterations, and rearrangements involving
268  CGC genes and genomic alterations associated with FDA approvals most frequently assigned a
269  nearest neighbor that shared drug sensitivity (19.7%, Figure 3b, Methods)[34].

270
271    This patient-to-cell line matchmaking module was then applied to our previously characterized
272    cohorts of patients with mCRPC and metastatic melanoma. Within the mCRPC cohort, the most
273    common nearest neighbor cell line among the 452 tested was VCaP, one of two prostate cancer
274    cell lines, for 25 of 150 patients. VCaP was sensitive to six therapies according to the GDSC;
275    however, these therapies (selisistat, SB52334, UNC0642, Trichostatin A, acetalax, and
276    linsitinib) do not have an established clinical role in mCRPC (Supplementary Figure 4). Nearest
277    neighbor cell lines to patients with metastatic melanoma were frequently sensitive to MEK and
278    RAF inhibitors, including SB590885 (BRAF inhibitor, nearest neighbor for 11 / 110 patients),
279    refametinib (MEK, 10), RAF_9304 (RAF, 8) and dabrafenib (BRAF, 7) (Figure 3c). Among
280    patients with metastatic melanoma that do not harbor *BRAF* p.V600E but do contain a *NRAS*
281    alteration (n = 24), the most common therapies which recurrent nearest neighbors were
282    sensitive to also included RAF_9304 (3 patients), refametinib (3), and SB590885 (3)
283    (Supplementary Figure 5).

284    **Integrated clinical interpretation in a prospective precision oncology trial**
285    We lastly compared therapeutic strategies nominated by the complete MOAlmanac
286    methodology with those administered to 83 patients in I-PREDICT (NCT02534675), a
287    prospective clinical trial evaluating personalized therapies based on panel sequencing
288    (Foundation Medicine's FoundationOne)[35]. Citations and relationships between molecular
289    features and clinical action from the study were reviewed and categorized by MOAlmanac
290    evidence levels (Supplementary Table 3). MOAlmanac processed the 524 molecular features
291    reported for I-PREDICT's 83 patients on a per-patient basis. Therapies administered in the
292    study (41 unique) or highlighted by our method (40) were categorized by therapeutic strategy
293    according to expert review based on shared pathway targets, resulting in a total of 31 unique
294    strategies (Supplementary Table 3). An overlap in recommended therapeutic strategy was
295    observed in 38 (46%) patients (Supplementary Figure 6). For patient therapy pairs highlighted
296    by MOAlmanac based on FDA evidence or clinical guidelines, 67% and 50%, respectively, were
297    involved in a therapeutic strategy administered by the study. Of the 13 patients with a therapy
298    highlighted by MOAlmanac associated with FDA approved or Guideline evidence that were not
299    involved in an overlapping strategy, 5 patients had another therapy which utilized a strategy
300    administered by I-PREDICT and the remaining 8 nominated therapies approved for other
301    disease contexts. For nominations based on weaker evidence categories, the concordance was
302    18% for preclinical and 50% for inferential (Figure 4a). The most common concordant strategies
303    were ER signaling inhibition, PI3K/AKT/mTOR inhibition, and immunotherapy (9, 9, and 7
304    patients, respectively). Of strategies that were not shared, I-PREDICT favored VEGF inhibition
305    for patients with *TP53* alterations (20 patients) whereas MOAlmanac frequently highlighted
306    assertions such as PRMT5 inhibition (13 patients) based on a preclinical relationship showing
307    efficacy of EPZ015666 for *CDKN2A/B* deletions (Figure 4b).
308
309    Finally, using our patient-to-cell line matchmaking module, nearest neighbor cell lines were
310    sensitive to a median of 2 therapies. For example, I-PREDICT administered everolimus and
311    MOAlmanac highlighted AZD8186 and pictilisib in the case of study id 105, a 60 year old female
312    with breast cancer. The nearest neighbor cell line, CAL-29 (bladder carcinoma), was sensitive to

313    taselisib and alpelisib as reported by GDSC2, both of which also target PI3K/Akt/mTOR. In
314    another case, I-PREDICT administered lenvatinib and ramucirumab for VEGF/VEGFR inhibition
315    to study id A009, a 44 year old male with esophageal adenocarcinoma. MOAlmanac highlighted
316    infigratinib for FGFR inhibition for therapeutic sensitivity and the nearest neighbor cancer cell
317    line, A204 (soft tissue), observes sensitivity to both VEGF and FGFR inhibition (VEGF:
318    cediranib, linifanib, motseanib, ponatinib, and tivozanib and FGFR: ponatinib). Thus,
319    MOAlmanac recapitulates established decision making paradigms in a prospective pan-cancer
320    setting and extends potential assertions in new therapeutic directions in other settings.

## 321    DISCUSSION

322    Here, we present a clinical interpretation method paired with a novel knowledgebase to facilitate
323    decision-making in precision oncology. In addition to first-order feature consideration,
324    MOAlmanac considers second-order molecular features such as mutational signatures, tumor
325    mutational burden, microsatellite stability, and ploidy, as well as high-throughput therapeutic
326    screens of cancer cell lines. Taken together, MOAlmanac addresses two key needs for
327    precision cancer medicine: 1) Point-of-care individualized patient treatment considerations
328    based on complex molecular interactions that considers evidence beyond FDA approvals and
329    clinical guidelines, and 2) Novel therapeutic hypotheses based on integrative interpretations that
330    can be evaluated in preclinical follow up and prospective trials. When applied to retrospective
331    cohorts, we observed that these novel features of MOAlmanac -- assessment of second-order
332    genomic features and consideration of preclinical or inferential evidence -- provided additional
333    hypotheses for prognosis and therapeutic sensitivity and resistance, especially for otherwise
334    variant-negative tumors.
335
336    While individual precision oncology studies require fixed versions of alteration-action knowledge
337    bases, rapidly expanding scope of literature on which these databases originate requires
338    constant updating that makes prospective assessment of precision oncology programs difficult.
339    This challenge was evident in comparing MOAlmanac to the I-PREDICT trial, as differences in
340    match selection were driven by differences in therapeutic availability at different time points,
341    variable knowledge capture of the vast precision oncology hypothesis landscape, and levels of
342    evidence to justify treatment selection. These results are suggestive of the urgency to
343    standardize genomic-based clinical trial data and aggregate knowledge bases to parse the vast
344    literature in precision oncology and enable principled, evidence-based clinical care[5,36]. Manual
345    curation of literature is inherently laborious, and prior efforts have encouraged crowdsourcing
346    and meta studies to address this challenge[4,5,37].
347
348    Furthermore, there were areas of note that could specifically improve our evaluation of patient-
349    to-cell line matchmaking for translational hypothesis generation. First, not all cell lines were
350    tested with every therapy; if they were, shared drug response could be characterized in a more
351    nuanced manner than the current boolean status. Second, there is likely an opportunity to
352    develop improved genomic similarity models which align with therapeutic sensitivity. The advent
353    of large, clinically annotated and molecular profiled patient cohorts may enable these
354    techniques and patient similarity networks to be evaluated for precision cancer medicine on

355 patient profiles rather than cancer cell lines[1,38,39]. Indeed, our primary motivation is to develop
356 similarity metrics that account for multiple data types from tumors to properly leverage nearest
357 neighbor approaches. These approaches, which prospectively leverage genomic data rather
358 than retrospectively curated data sources, are imperative to develop therapeutic hypotheses for
359 patients who are variant negative.
360
361 In conclusion, MOAlmanac catalyzes the use of expanded feature types, evidence sources, and
362 algorithms for clinical interpretation of integrative molecular features for precision cancer
363 medicine applications. Incorporation of MOAlmanac into future translational studies and clinical
364 trials may directly enable evaluation of the precision oncology hypothesis across patient
365 populations. Furthermore, MOAlmanac can promote evaluation of patient similarity networks
366 using both clinical and preclinical knowledge to aid precision cancer medicine at the individual
367 patient level for translational discovery. The Molecular Oncology Almanac is available at
368 https://moalmanac.org. This method is available on Github
369 (https://github.com/vanallenlab/moalmanac), Docker Hub
370 (https://hub.docker.com/r/vanallenlab/moalmanac), and on the Broad Institute's Terra
371 (https://portal.firecloud.org/#methods/vanallenlab/moalmanac/). In addition, a web portal to
372 process individual cases through a user interface atop of Terra is available at
373 https://portal.moalmanac.org/. All code related to analyses and figures herein can be found on
374 Github (https://github.com/vanallenlab/moalmanac-paper). Finally, to facilitate crowdsourced
375 updating of MOAlmanac's knowledge base, Molecular Oncology Almanac Connector (a Google
376 Chrome extension) is available to enable users to nominate relationships with minimal effort.

# METHODS

**Creating a knowledge base**

*Defining a database schema*

380 An SQL schema was planned and abstracted with Vertabelo for cataloging clinical assertions
381 relating molecular features to clinical action. The schema contained four primary abstractions:
382 Assertion, Feature, Source, and Version with additional tables to relate assertion to features and
383 sources; Assertion_To_Feature and Assertion_To_Source, respectively (Supplementary Figure
384 7). The underlying data structure is implemented as an SQLite database and managed with
385 Python and SQL Alchemy.
386
387 The Assertion table is used to catalog a given clinical action. The context of an assertion is
388 catalogued with disease as described in the source (disease), which is mapped to an oncotree
389 code (oncotree_code) and term (oncotree_term), and any applicable disease context such as
390 disease stage (context). If regarding therapeutic sensitivity or therapeutic resistance, the drug
391 name is entered (therapy_name) along with its type (therapy_type: targeted therapy,
392 chemotherapy, radiation, immunotherapy, hormonal therapy, or combination) and a boolean
393 integer of 1 for asserting a relationship to differential therapeutic sensitivity or resistance or 0
394 for asserting no such relationship. This data structure allows MOAlmanac to capture negative

395 studies documenting that a given feature is not associated with differential therapeutic
396 sensitivity). If regarding prognosis, a boolean integer is entered to suggest a favorable or
397 unfavorable prognosis (favorable_prognosis). The evidence of an assertion is recorded
398 (predictive_implication); available values are "FDA-approved", "Guideline" for clinical guideline,
399 "Clinical trial" for associations reported from clinical trials, "Clinical evidence" for retrospective
400 studies or human studies not directly reported from a clinical trial, "Preclinical evidence" for
401 findings from mouse models or cancer cell lines, or "Inferential evidence" for findings from
402 mathematical models or an association between molecular features. In some cases, we denote
403 favored assertions (preferred_assertion) to "tie break" otherwise equal assertions based on
404 published literature and clinical use; e.g. Dabrafenib and Trametinib over Vemurafenib for *BRAF*
405 p.V600E. A free text description of the clinical assertion is curated for all entries (description)
406 along with an entry date (created_on) and last modified date (last_updated).
407
408 Molecular features are associated with assertions and are catalogued in a flexible manner to
409 accommodate different attributes of a feature type using feature definitions. For example,
410 rearrangements are defined as having a rearrangement type (translocation, fusion), participating
411 genes (gene1, gene2), and a locus; separately, copy number alterations are defined as having a
412 gene, direction, and cytoband. Rearrangements, somatic variants, germline variants, copy
413 number alterations, microsatellite stability, mutational signatures, mutational burden, neoantigen
414 burden, knockdown, silencing, and aneuploidy are currently catalogued with feature definitions.
415 New feature definitions may be easily programmatically defined, allowing the rapid addition of
416 new features without having to modify the underlying data schema.
417
418 Sources are catalogued such that all sources will be associated with a citation, source type
419 (abstract, FDA, guideline, journal), and url. Journal articles are further annotated with the
420 associated PubMed ID (PMID) and DOI. Sources regarding a clinical trial will catalog the
421 National Clinical Trial (NCT) registry number.
422
423 Version is an unconnected table used to catalog major, minor, and patch numbers of the
424 database.


425 *Iterating from TARGET*

426 TARGET catalogued clinical assertions primarily by gene associated with types of recurrent
427 alterations and examples of therapeutic agents paired with an aggregate rationale for the gene.
428 Literature review was performed by curators to review FDA approvals, clinical guidelines, and
429 journal articles to associate clinical assertions from TARGET with a citation. Associations to 52
430 genes were removed due to insufficient evidence, recent evidence conflicted with the underlying
431 assertion for 1 gene, and 5 genes were partially retained. Ten genes were not migrated to
432 MOAlmanac because we chose to not catalog the underlying assertion type; specifically, we
433 intentionally chose to not include diagnostic relationships and we reclassified biallelic loss to
434 copy number deletions.

435  *Cataloging additional assertions*

436  Subsequent curation efforts cataloged FDA approvals, clinical guidelines, conference abstracts,
437  or recently published literature. Relationships were further categorized by the clinical implication
438  of the assertion (therapeutic sensitivity or resistance or prognostic value), therapy type if
439  relevant, and evidence. Genomic feature types considered were somatic and germline variants,
440  copy number alterations, rearrangements, mutational burden, COSMIC mutational signatures,
441  microsatellite stability status, and aneuploidy.
442
443  The knowledge base contained 722 assertions which relate molecular features to therapeutic
444  response and prognosis and 4 related to adverse event risk, manually curated from literature
445  review of FDA approvals (87 assertions), clinical guidelines (187), published journal articles
446  (446), and abstracts (5). In addition to characterizing targeted therapies (417 relationships), we
447  have catalogued relationships related to immunotherapies (48), chemotherapies (40), radiation
448  (19), hormonal treatments (7), and combination therapies (11, Figure 1c).
449
450  No further assertions were added to MOAlmanac past March 23rd, 2020 for the purposes of this
451  study.


452  *Comparison to other knowledge bases*

453  Molecular Oncology Almanac was categorically compared to CIViC and OncoKB, two similar
454  precision oncology knowledge bases, across the categories of therapy types, molecular feature
455  types, assertion types, catalogued evidence, curation type, accessibility, number of assertions,
456  and counted therapy types (Supplementary Table 4**)**. Citations with PubMed reference numbers
457  (PMIDs, 458 citations) were compared and we observed similar findings to previous meta-
458  studies, that no one database subsumes another (Supplementary Figure 8)[37].


459  **Developing a clinical interpretation method**

460  *Accepted inputs*

461  Molecular Oncology Almanac accepts any combination of somatic variants, copy number
462  alterations, rearrangements, germline variants, somatic variants from another source such as a
463  validation sequencing, and breadth of coverage. In addition, several single value or boolean
464  features are passable such as the purity and ploidy of the tumor as float values, a categorical
465  input for microsatellite stability status, a boolean flag to note whole genome doubling. Free text
466  fields are also available to enter a patient or sample id, tumor type, stage, and general
467  description of the molecular profile.
468
469  Input files to MOAlmanac have expectations on their format, which can be found on the
470  method's Github. Somatic, both primary or validation sequencing, and germline variants
471  conform to the National Cancer Institute's Genomic Data Commons MAF v1.0.0 format,
472  requiring: Hugo_Symbol, Chromosome, Start_position, End_position, Reference_Allele,
473  Tumor_Seq_Allele1, Tumor_Seq_Allele2, Variant_Classification, Protein_Change,
474  Tumor_Sample_Barcode, Normal_Sample_Barcode, t_ref_count, t_alt_count. MOAlmanac is

475 coded to accept input columns based on Oncotator for these inputs; however, this can be
476 changed by editing the colnames.ini file[40]. MOAlmanac currently is coded to accept total copy
477 number alterations produced by ReCapSeg, or GATK3 CNV, and annotated by Oncotator,
478 requiring the columns gene, segment_contig, segment_start, segment_end, sample, and
479 segment_mean[41]. MOAlmanac is coded to accept rearrangements directly from STAR Fusion,
480 requiring the columns fusion_name, SpanningFrags, LeftBreakPoint, and RightBreakPoint[42].
481 Breadth of coverage is the sum of calculable bases used to call somatic variants, and is
482 required to calculate nonsynonymous mutational burden; a text file containing the integer, such
483 as summing MuTect 1.0's call stats output, suffices.
484
485 The input arguments stage, purity, ploidy, and description are only used for display as metadata
486 in the produced actionability report. Provided tumor types are mapped to standardized ontology
487 terms and codes using Oncotree (http://oncotree.mskcc.org/#/home), if possible. Patient ID is
488 also used as metadata and is also used as a prefix to label all generated outputs.


489 *Annotation and evaluation of individual molecular features*

490 Somatic variants, copy number alterations, and gene fusions are annotated with MOAlmanac,
491 Cancer Hotspots, 3D Hotspots, Cancer Gene Census (CGC), Molecular Signatures Database
492 (MSigDB), and COSMIC [18,19,21–23]. Genomic events are first annotated for their gene presence (1
493 for present, 0 for wild type) and then receives a higher integer score if applicable; for example,
494 somatic variants whose protein change appears in Cancer Hotspots will be noted by a 2.
495 Somatic and germline variants are also annotated with ClinVar and ExAC to identify pathogenic
496 or likely pathogenic variants and common variants [24,25]. Somatic variants and copy number
497 alterations are annotated and evaluated based on a heuristic similar to PHIAL, sorting to events
498 based on their presence in data sources (Supplementary Figure 1a).
499
500 MOAlmanac considers individual non-synonymous variants (missense, nonsense, nonstop,
501 frameshift, insertions, and deletions), copy number alterations that are outside of 1.96 standard
502 deviations from the mean of unique segment means (above 97.5 percentile for amplifications
503 and below 2.5 percentile for deletions), and at least 5 spanning fragments for fusions. Events
504 which meet these criteria will be scored by MOAlmanac's somatic heuristic and be provided in
505 the output file with the suffix ".somatic.scored.txt", while filtered alterations are made available in
506 the output noted by the ".somatic.filtered.txt" suffix.
507
508 For genomic alterations whose gene appears in Molecular Oncology Almanac, the clinical
509 relevance will be labeled based on the match to the catalogued molecular feature and evidence
510 tier of the matched relationship. Complete matches to explicit features (e.g. protein change for
511 variants, direction for copy number alteration, or fusion and partner) will be labeled as Putatively
512 Actionable whereas partial matches or incompletely characterized features (the gene is
513 catalogued of that data type; e.g. a *ETV6-NTRK1* fusion matches to an assertion of *NTRK1*
514 fusions) is labeled as Investigate Actionability. If an alteration's gene appears in Molecular
515 Oncology Almanac but not catalogued as the same data type, the alteration will be labeled as
516 Biologically Relevant and is not associated with any clinical relationships. For each provided
517 genomic feature, a match is searched for relationships associated with therapeutic sensitivity,

518 resistance, and disease prognosis and, if either labeled as Putatively Actionable or Investigate
519 Actionability, evidence level of the association, therapy name and therapy type (if sensitivity or
520 resistance) or favorable prognosis, relationship description, citation, and URL for the citation are
521 associated. These actionable features are made available in the output file with the suffix
522 ".actionable.txt".

523

524 In addition, a few outputs regarding germline variants are highlighted and made available, if
525 provided (Supplementary Figure 1b). Variants in genes related to hereditary cancers, based on
526 a panel of 83 genes commonly used for germline testing, are produced in an output with the
527 suffix ".germline.hereditary_cancers.txt"[43]. Likewise, variants in genes noted by the American
528 College of Medical Genetics and Genomics secondary findings v2 [44] are highlighted in the
529 output with the suffix ".germline.acmg.txt". Lastly, germline variants in genes related to somatic
530 cancers (based on a gene presence in MOAlmanac, Cancer Hotspots, or Cancer Gene Census)
531 are noted in the output of the suffix ".germline.cancer_related.txt". Germline variants which
532 match to MOAlmanac will also be included in the actionable output if (1) they are not labeled as
533 common in ExAC (an allele frequency greater than 1 in 1,000 alleles), (2) are labeled as a
534 pathogenic or likely pathogenic variant in ClinVar, or (3) a truncating (frameshift, nonsense,
535 nonstop, or splice site) variant.

536

537 If somatic single nucleotide variants are provided for both primary and secondary (also referred
538 to as validation or orthogonal sequencing) sequencing, MOAlmanac will annotate variants called
539 in the primary sequencing based on their presence (allelic fraction and coverage) in the
540 secondary sequencing. The power to detect variants in the secondary sequencing is calculated
541 using a beta-binomial distribution with $k$ equal to 3 for a minimum of three reads, $n$ as coverage
542 of the variant in secondary sequencing, $alpha$ and $beta$ defined as the alternate and reference
543 read counts + 1 as observed from the primary sequencing, respectively. This approach is
544 consistent with best practices by Yizhak et al. 2019 with RNA MuTect[11]. The allelic fraction of
545 somatic variants observed in primary and orthogonal sequencing are plotted against each other
546 in a scatter plot in the output of the suffix ".validation_overlap.png", with variants observed with
547 detection power greater than or equal to the specified minimum (default 0.80) colored in blue
548 and those otherwise grey. At the moment, MOAlmanac only leverages orthogonal sequencing
549 for validation and does not use it for discovery. When applied to the retrospective cohorts of
550 metastatic melanoma and mCRPC, we had sufficient power to observe 190 of 453 applicable
551 clinically relevant variants. Of note, *AR* p.L702H and p.T878A, variants putatively associated
552 with resistance to androgen deprivation, were observed in the RNA of 6 and 4 patients,
553 respectively[45].


554 *Annotation and evaluation of integrative and second-order genomic features*

555 To ease the process of reviewing multiple intra-gene alterations, MOAlmanac summarizes all
556 somatic variants, germline variants, copy number alterations, and fusion events per gene for
557 genes found within MOAlmanac, Cancer Hotspots, and Cancer Gene Census. Any genes with
558 at least one alteration across any data type will be reported in the output with the suffix
559 ".integrated.summary.txt".
560

561 Somatic alterations are annotated with the number of frameshift, nonstop, nonsense, or splice
562 site germline events within the same gene. This count is labeled as the column
563 "number_germline_mutations_in_gene" in the output of the suffix ".somatic.scored.txt".
564
565 Tumor mutational burden (TMB) is calculated based on the number of nonsynonymous variants
566 divided by the somatic calculable bases. TMB is compared to values calculated for TCGA
567 molecular profiles by Lawrence et al. 2013 to yield a pancan percentile and tissue-specific
568 percentile, if ontology matched to one of the 27 tumor types studied in the publication[46]. TMB for
569 a molecular profile is designated as high if greater than 10 nonsynonymous variants per
570 megabase and greater than or equal to the 80th tissue-specific percentile, or pancan percentile
571 if not mapped.
572
573 COSMIC mutational signatures are evaluated using deconstructSigs by running R as a
574 subprocess using the default trinucleotide counts method [47,48]. Signatures with a contribution
575 greater than a specified minimum contribution (default: 0.20) are annotated at least as
576 Biologically Relevant and annotated using MOAlmanac for consideration of actionability.
577 Nucleotide context counts are made available in table format directly from deconstructSigs as
578 an output with the suffix ".sigs.context.txt" and signature contributions with the suffix
579 ".sigs.cosmic.txt". Trinucleotide counts of a considered molecular profile are plotted based on
580 raw and normalized counts in the outputs ".sigs.tricontext.counts.png" and
581 ".sigs.tricontext.normalized.png", respectively.
582
583 Microsatellite stability is both directly considered as a categorical input for status and indirectly
584 by highlighting potentially related variants. As a direct input, users may flag microsatellite status
585 as microsatellite stable, microsatellite instability low, microsatellite instability high, or unknown.
586 Genomic alterations which appear in genes related to microsatellite instability are highlighted as
587 supporting variants and Biologically Relevant and further noted in their own output, with the
588 suffix ".msi_variants.txt"; specifically, the genes considered are *ACVR2A*, *DOCK3*, *ESRP1*,
589 *JAK1*, *MLH1*, *MSH2*, *MSH3*, *MSH6*, *PMS2*, *POLE*, *POLE2*, *PRMD2*, and *RNF43*[49,50]. As of this
590 publication, MOAlmanac has only catalogued assertions related to MSI-High status.
591
592 Whole genome doubling, or aneuploidy, is available for consideration as a boolean-valued input
593 and, if flagged, will evaluate for clinical relevance based on the currently catalogued assertions.
594 As of this publication, MOAlmanac has catalogued Bielski et al. 2018's observation that whole
595 genome doubling being associated with adverse survival across a pan-cancer setting[28].
596
597 Mutational burden, mutational signatures, microsatellite stability, and whole genome doubling
598 are at most highlighted as Investigate Actionability by Molecular Oncology Almanac for clinical
599 assessment.


600 *Creating clinical actionability reports*

601 Clinical actionability reports are created for all profiles processed with Molecular Oncology
602 Almanac, generated with Python 3.6, Flask, and Frozen Flask.
603

604 The reports contain sections containing profile metadata (Profile Information), molecular
605 features associated as a Putatively Actionable or Investigate Actionability predictive implication
606 for therapeutic sensitivity or resistance and prognosis, as well as variants associated with
607 Biological Relevance (Actionability Report). Associations list the implication, evidence, and
608 associated therapy and description of clinical assertion as rationale. Sources for each
609 association are available as hyperlinks labeled as "[source]", equivalent assertions are available
610 to view in a modal labeled , and preclinical efficacy of the assertion is also available as modal, if
611 applicable.
612
613 The 5 most similar cell lines to the provided molecular profile are listed by their CCLE name
614 along with their sensitive therapies and clinically relevant features. For each cell line, a modal is
615 available that lists their Broad/DepMap and Sanger Institute aliases and somatic variants, copy
616 number alterations, and fusions in any MOAlmanac, Cancer Hotspot, or CGC gene as well as
617 the ln(ic50), AUC, and z score for each of the top 10 most sensitive therapies of the cell lines.
618 This feature can be hidden in the clinical report passing diable_matchmaking as a parameter to
619 the method.
620
621 Due to being produced with Frozen Flask, these web based reports are a single html file with no
622 additional file dependencies. They usually are no larger than 1 Mb in size.


### Comparing PHIAL-TARGET and MOAlmanac with two retrospective studies

*Data acquisition and sample processing*

625 Whole-exome sequencing (WES) and RNA sequencing (RNA-seq) was acquired for 110
626 previously published patients with metastatic melanomas (n = 44 with RNA)[26] and 150 patients
627 with castration-resistant prostate cancers (mCRPC, n = 149 with RNA)[27]. Subsequent sample
628 processing was performed on the Broad Institute and Verily Life Sciences' Terra Google Cloud
629 platform.
630
631 Whole-exome sequencing was used to call somatic and germline variants and copy number
632 alterations. MuTect 1.0 was used to identify single nucleotide variants (SNVs) and somatic
633 calculable bases of individual tumor samples while Strelka was used to identify insertions and
634 deletions (InDels)[51,52], run utilizing the Getz Lab CGA WES Characterization pipeline at the
635 Broad Institute. Artifacts introduced by DNA oxidation during the sequencing process were
636 removed [53]. Mutations calls were compared to a panel of germline samples and were removed if
637 they appeared in more than three germline samples[54]. Germline variants were called using
638 Deep Variant[55]. Segmented total copy number was calculated across the exome by comparing
639 fractional exome coverage to a panel of normals using CapSeg [56,57]. Tumor purity and ploidy
640 was calculated using FACETS[58].
641
642 Transcriptome BAMs were converted to FASTQ format and aligned using STAR [59]. Fusions
643 were then called using STAR Fusion[60]. STAR aligned bams were calibrated following GATK's
644 best practices for variant discovery in RNA-seq (https://github.com/broadinstitute/gatk-
645 docs/blob/3333b5aacfd3c48a87b60047395e1febc98c21f9/gatk3-methods-and-

646   algorithms/Calling_variants_in_RNAseq.md) using GATK 3.7[61–63]. Somatic variants observed in
647   whole-exome data were then force called from the recalibrated RNA-seq bams for each
648   individual using MuTect 1.0.
649
650   Somatic variants from both WES and RNA-seq, germline variants, and copy number alterations
651   were annotated using Oncotator v1.9.1[40].

652   *Comparison of clinically relevant events*

653   Molecular features were processed for all 260 samples by both PHIAL 1.0.0
654   (https://github.com/vanallenlab/phial)[2] and MOAlmanac. While both methodologies considered
655   all available genomic events, PHIAL considered somatic variants and copy number alterations
656   while MOAlmanac additionally considered germline variants, rearrangements, mutational
657   burden, mutational signatures, and whole-genome doubling. Microsatellite stability was not
658   considered for this analysis as labels from testing, if performed, were not available. Events that
659   matched with the underlying knowledge base as either Investigate Actionability or Putatively
660   Actionable, thus stronger than simply a gene match, were considered for clinical relevance
661   (Supplementary Figure 2). While the differences were impacted by literature curation and
662   MOAlmanac considering additional feature types, they were also impacted by changing how
663   copy number alterations are handled; PHIAL called copy number alterations based on a
664   threshold, |segment mean| $\geq$ 1, whereas MOAlmanac utilizes a percentile approach, top or
665   bottom 2.5%.

666   **Expanded methods for directly leveraging preclinical models**

667   *Data acquisition and processing*

668   Somatic variants and copy number alterations for cancer cell lines catalogued in the Cancer Cell
669   Line Encyclopedia were gathered from cBioPortal and fusions and therapeutic sensitivity were
670   downloaded from the Sanger Institute's Genomics of Drug Sensitivity in Cancer (GDSC) [32,33].
671   Cancer cell lines were standardized by name and filtered for by requiring: all four data types
672   being available, being of solid tumor origin, not subject to genetic drift between Broad and
673   Sanger versions of the cell line per Ghandi et al. 2019, and not reclassified as fibroblast like by
674   Weck et al. 2017 and Ghandi et al. 2019 [32,64]; resulting in 452 cancer cell lines. Somatic
675   variants, copy number alterations, and fusions were formatted for usage and annotated by
676   Molecular Oncology Almanac.

677   *Directly leveraging preclinical models to evaluate efficacy*

678   All GDSC1 and GDSC2 therapies were mapped to therapies catalogued in MOAlmanac. For all
679   therapies associated with genomic events by MOAlmanac for which a GDSC mapping exists, a
680   sensitivity dictionary is created in which each key is associated with a clinically relevant feature
681   found by the method. For each feature, we list all mutant and wild type cell lines for each
682   component; e.g. when considering *CDKN2A* deletions, mutant and wild type lists are made for
683   all cell lines that have any alteration in *CDKN2A* (somatic variant, copy number alteration, or
684   fusion), cell lines that have a *CDKN2A* copy number alteration, and cell lines that have a

685      *CDKN2A* deletion. For each pairing of mutant and wild type cell lines, the IC50 values are
686      compared with a Mann-Whitney-Wilcoxon test to evaluate if a significant difference exists
687      between the two distributions. A box plot of mutant and wild type cell lines and their IC50 values
688      is also created, labeled by the genomic feature used to stratify.
689
690      The results of such testing are reported in two outputs, the actionability report with the suffix
691      ".report.html" and a table compiling all examinations with the suffix ".preclinical.efficacy.txt".
692      When applicable, a hyperlink labeled as "[Preclinical evidence]" will appear under "Therapy &
693      rationale" for variants and features associated with therapeutic sensitivity. Upon clicking the link,
694      a modal window opens showing all box plots of comparisons along with the number of wild type
695      cell lines, number of mutant cell lines, and the Mann-Whitney-Wilcoxon statistic and p-value for
696      each feature evaluated. In addition, IC50 median, mean, and standard deviation can be found
697      for all relationships evaluated in the mentioned preclinical efficacy table output.

698      *Directly leveraging preclinical models for patient-model matchmaking*

699      We sought to directly leverage molecular profiles for clinical interpretation. For the purposes of
700      this application, we sought to compare a case molecular profile to a larger population and sort
701      other members by genomic features such that the nearest neighbor to our case profile shared
702      drug sensitivity. In absence of a large cohort of clinically annotated primary or metastatic tumor
703      profiles, we utilized cancer cell lines which have been characterized by high throuput drug
704      screens and evaluated by comparing cell lines against cell lines.
705
706      GDSC z scores of therapies applied to cell lines were utilized to convert continuous valued IC50
707      response curves to boolean valued sensitive (z score $\leq$ -2) or resistant (z score $\geq$ 2)[33]. Pairwise
708      comparisons were made between all cell lines which contained GDSC therapeutic response
709      data, noting the intersection of therapies which both profiles were deemed sensitive to as well
710      as the intersection size. If the intersection size was greater than 0, the pair was deemed to
711      share therapeutic sensitivity. When evaluating a novel case profile the matchmaking module of
712      MOAlmanac, the 452 cancer cell lines, which result from filtering described in two sections prior,
713      are used for comparison. However, for evaluation, we further required that cell lines are
714      sensitive to at least one therapy and that there exists at least one other cell line that shares
715      therapeutic sensitivity, so that there is at least one true positive when sorting other cell lines,
716      resulting in 377 cell lines.
717
718      After somatic variants, copy number alterations, and fusions were annotated and evaluated by
719      MOAlmanac, molecular features were vectorized into sample x feature tables. The coding of
720      features was dependent on the model implemented, discussed more explicitly in the next
721      section; however, some commonalities exist. All elements were boolean valued and thus all
722      feature tables were sparse boolean arrays. When a similarity model involved genes, either the
723      CGC (n = 719 genes) or MOAlmanac (130) gene sets were used. Among a series of models
724      tested, we found the best performing model to be using Similarity Network Fusion on four
725      sample x feature tables: CGC genes altered by somatic variants, copy number alterations, and
726      fusions and a fourth table of samples x specific molecular features associated with an FDA
727      approved therapy, subsequently referred to as SNF: CGC & FDA.

728
729     Evaluation metrics were borrowed from ranked retrieval.
730
731     The performance of how a similarity metric sorts cell lines relative to one cell line are evaluated
732     using precision @ rank ($k$), recall @ $k$, and average precision. Consider four cell lines sorted in
733     order relative to a case profile such that the first and third share therapeutic sensitivity with the
734     case profile and the second and fourth does not (Figure 3a). Cell lines which share therapeutic
735     sensitivity can be considered relevant. To calculate precision @ $k$, given $k$ neighbors, we divide
736     the number of relevant neighbors divided by $k$; e.g. considering the first neighbor (k=1) yields a
737     precision @ 1 of 1.0 (1 relevant neighbor / 1) but considering the second neighbor as well yields
738     a precision @ 2 of 0.5 (1 relevant neighbor / 2). Recall is calculated as the fraction of overall
739     relevant neighbors returned when considering $k$ neighbors; at $k = 1$ recall is calculated to be 0.5
740     in our example, until $k = 3$ when a second relevant cell line is returned thus recall is calculated
741     to be 1.0, and recall = 1.0 at $k = 4$. Average precision (AP) is calculated by taking the average of
742     precision values at positions of a relevant neighbor; using our example, relevant neighbors exist
743     at precision @ k = 1 and 3 with associated precision values of 1.0 and 0.66 so the average
744     precision for this sort, or query to use terminology from information retrieval, is calculated to be
745     0.83.
746
747     The performance of a similarity metric for many queries can be evaluated by calculating the
748     mean average precision (mAP). Given three case profiles which sorted cell lines against them
749     with average precision values of 0.66, 0.565, and 0.25, the mean average precision is the
750     average of them, which is calculated to be 0.492. In our context, for each similarity model, we
751     calculate the average precision for each cell line and the mean average precision across all cell
752     lines (Supplementary Table 5).
753
754     Models can be compared pairwise with permutation testing (Supplementary Table 6). The
755     difference in mean average precision (delta mAP) is chosen as a test statistic and the AP @ k
756     values are shuffled for all 377 values of k. Given these shuffled AP @ k values, mAP values are
757     calculated along with a delta mAP and the delta mAP is recorded. This was performed over
758     10,000 iterations using seeds 0 to 9,999 to create a distribution of delta mAP values. The test
759     statistic is compared to the distribution to generate a p-value and, if the p-value was $\geq$ 0.05, it
760     was deemed that the two models were within the noise range of one another. Our best
761     performing model SNF: CGC & FDA was within the noise range of two other models, a multi-
762     pass sort of first using agreement based measure of molecular features associated with an FDA
763     approved therapy followed by agreement based sort of CGC genes mutated by any feature type
764     (Multi-pass sort: FDA & CGC, p=0.4013) and sorting cell lines by their mutant and wild type
765     status of variants in order based on the somatic heuristic in MOAlmanac (Somatic tree,
766     p=0.5458); however, SNF: CGC & FDA observed a stronger AP @ k = 1 in both cases, 0.193
767     versus 0.164 and 0.119, respectively.
768
769     There are several areas which we note that this framework could be improved. First, not all cell
770     lines were treated with all therapies and we can not deem an untested pair as sensitive or not
771     sensitive unless we resort to estimating missing data, thus, we assume that cell lines do not

772    respond to therapies which they were not tested to be conservative in our analysis. In the
773    setting of a complete pairing (all cell lines are treated with all therapies) we could incorporate a
774    more nuanced label. For example, we could continue using the z score thresholds but instead
775    label based on the jaccard index of shared therapies or we could transition to using a
776    continuous valued similarity of drug sensitivity such as euclidean distance of IC50s or perform a
777    PCA. In either case, a complete pairing of therapies and cell lines would enable us to use
778    additional evaluation metrics such as Discounted Cumulative Gain (DCG), ranking other cell
779    lines based on a relevance scale rather than a boolean condition and rank. Secondly, rather
780    than evaluate cell lines against cell lines, we envision that an ideal experiment for this analysis
781    would involve a cohort of paired primary tumor samples and patient derived cell lines in which
782    we would hope that the paired patient derived cell line would be deemed most similar to its
783    corresponding tissue sample. Such a setting would enable the studying of performance as a
784    function of cell line passages. Expression was not used in this analysis as it is a feature
785    modality not yet commonly used at the point-of-care.

786    *Models and calculating similarity metrics*

787    Several models were implemented to characterize similarity between cancer cell lines based on
788    genomic features. Models were evaluated using average precision, specifically average
789    precision @ k = 1, and mean average precision. In short, our best performing model (SNF: FDA
790    & CGC) observed a AP @ k = 1 of 0.194 which was 2.03x better than random but still only
791    recommends a nearest neighbor for one fifth of cell lines. We are excited to see improvements
792    in directly leveraging molecular profiles for clinical interpretation. Performance of models can be
793    found in Supplementary Table 5. Models include, listed alphabetically:
794
795    Compatibility (compatibility). Inspired by dating algorithms, we weigh each molecular feature (or
796    question) based on strength of the match (e.g. a BRAF deletion only matches BRAF p.V600E
797    by gene). With these relative weights, we calculate a max score for each sample and compare
798    against other cell lines.
799
800    Jaccard of MOAlmanac feature types (jaccard-almanac-feature-types). We sort by agreement
801    based measure (jaccard) by considering both gene and data type for all somatic variants, copy
802    number alterations, and rearrangements catalogued in the Molecular Oncology Almanac (e.g.
803    CDKN2A copy number alterations match but not a CDKN2A deletion and CDKN2A nonsense
804    somatic variant).
805
806    Jaccard of MOAlmanac features (jaccard-almanac-features). We sort by agreement based
807    measure (jaccard) by considering all somatic variant, copy number, and rearrangement
808    molecular features catalogued in the Molecular Oncology Almanac.
809
810    Jaccard of MOAlmanac genes (jaccard-almanac-genes). We sort by agreement based measure
811    (jaccard) by considering any somatic variant, copy number alteration, and rearrangement in any
812    gene catalogued in Molecular Oncology Almanac.
813

814 Jaccard of CGC feature types (jaccard-cgc-feature-types). We sort by agreement based
815 measure (jaccard) by considering variants in a Cancer Gene Census gene and feature type
816 (e.g. CDKN2A copy number alterations match but not a CDKN2A deletion and CDKN2A
817 nonsense somatic variant).
818
819 Jaccard of CGC genes (jaccard-cgc-genes). We sort by agreement based measure (jaccard) by
820 considering any variant in a Cancer Gene Census gene.
821
822 Multi-pass sort: FDA & CGC (multi-pass-sort_fda-cgc). A weakness of agreement based
823 measure is that there will be tied values. We tie break similarities based on Molecular Oncology
824 Almanac features associated with FDA evidence by using similarity based on CGC genes.
825
826 Nonsynonymous variant count (nonsynonymous-variant-count). We assign neighbors based on
827 the absolute value of the difference of the number of coding somatic variants. This is a proxy for
828 mutational burden, because we do not have the number of somatic bases considered when
829 calling variants to use a denominator.
830
831 PCA of MOAlmanac genes (pca-almanac-genes). We run PCA and then nearest neighbors for
832 the vectorization of MOAlmanac genes, with mutants being without consideration of feature
833 type. For example, there is one feature called "TP53" and both TP53 nonsense variants and
834 copy number deletions can populate the element.
835
836 PCA of CGC genes (pca-cgc-genes). We run PCA and then nearest neighbors for the
837 vectorization of CGC genes, with mutants being without consideration of feature type. For
838 example, there is one feature called "TP53" and both TP53 nonsense variants and copy number
839 deletions can populate the element.
840
841 Random (random_mean). Randomly shuffle cell lines against one another across 100,000
842 seeds. This uses the seed of the average mean average precision.
843
844 SNF: MOAlmanac (snf_almanac). Rather than collapse all data types into a single similarity
845 matrix (e.g. with columns such as CDKN2A somatic variant, CDKN2A copy number alteration),
846 we use the python implementation of Similarity Network Fusion by Ross
847 Markello(https://github.com/rmarkello/snfpy)[34]. We fuse networks that describe agreement
848 based on variants in almanac genes in (1) somatic variants, (2) copy number alterations, and (3)
849 rearrangements.
850
851 SNF: CGC (snf_cgc). Rather than collapse all data types into a single similarity matrix (e.g. with
852 columns such as CDKN2A somatic variant, CDKN2A copy number alteration), we use the
853 python implementation of Similarity Network Fusion by Ross Markello
854 (https://github.com/rmarkello/snfpy)[34]. We fuse networks that describe agreement based on
855 variants in CGC genes in (1) somatic variants, (2) copy number alterations, and (3)
856 rearrangements.
857

858    SNF: FDA & CGC (snf_fda-cgc). We perform similarity network fusion using the python
859    implementation by Ross Markello (https://github.com/rmarkello/snfpy) to fuse networks that
860    contain: (1) CGC genes that contain a somatic variant, (2) CGC genes that contain a copy
861    number alteration, (3) CGC genes that contain a rearrangement, (4) Almanac features
862    associated with FDA evidence[34].

864    SNF: FDA & CGC genes (snf_fda-cgc-genes). We perform similarity network fusion using the
865    python implementation by Ross Markello (https://github.com/rmarkello/snfpy) to fuse networks
866    that contain (1) almanac features associated with FDA evidence and (2) any variant occurring in
867    a Cancer Gene Census gene.

869    Somatic tree (somatic-tree). This is somewhat inspired by CELLector by Najgebauer et al.[16].
870    One issue with agreement based measures is that each feature is weighted the same.
871    CELLector has a sorted list of genes/variants based on cancer type and will report similar cell
872    lines based on mutant / wild type status of each gene. While not exactly the same, we use the
873    annotations from various data sources appended to variants by Molecular Oncology Almanac to
874    create a priority list for variants (hotspots ranked the highest, etc.). For each case sample, we
875    consider the genes which are observed to be mutated and preserve the order that they would
876    appear in the somatic.scored.txt output of MOAlmanac. All other samples are then sorted by
877    their mutant / wild type status of these genes.


878    **Comparing to a prospective clinical trial, I-PREDICT**

879    We compared the clinical actions administered based on molecular profiles to patients in the I-
880    PREDICT prospective clinical trial to those highlighted by Molecular Oncology Almanac[35]. All
881    genomic events considered were present in the supplementary text of the study and we
882    extracted molecular features, therapies administered, and citations. Disease ontologies were
883    mapped to Oncotree terms and codes (http://oncotree.mskcc.org/). Molecular features were
884    formatted for annotation and evaluation by MOAlmanac.

886    Citations providing rationale for therapies administered based on molecular features were
887    extracted from the supplementary text, obtained, read, commented on, and categorized by
888    evidence level. Molecular features considered by the study were merged with annotations made
889    by MOAlmanac and, using the author notes from the supplementary text, we annotated if the
890    study targeted the molecular feature. Therapy and associated molecular features were mapped
891    to therapeutic strategies by expert review. Therapies administered in the study and those
892    highlighted by MOAlmanac for therapeutic sensitivity were listed on a per patient basis and
893    evidence levels were annotated for each therapy per patient. For therapies administered by the
894    study, citations cited per patient were referenced again for the specific relationship between
895    therapeutic strategy or therapy and molecular feature. Each therapy administered was binned
896    based on the evidence level or annotation as no citation, if the therapy was administered not on
897    the basis of molecular features, or citation listed not applicable, if the citation(s) listed did not
898    mention the therapy, strategy, or target. In some cases which would have resulted in the latter,
899    we transcribed that perhaps a source cited for another relationship in the cohort and cited that
900    source. Therapies were tagged with a boolean value if they were involved in a shared

901  therapeutic strategy between what was administered in I-PREDICT and highlighted by
902  Molecular Oncology Almanac for a given patient (Supplementary Table 3).

903  **Web-based tools to improve accessibility**

904  *Browsing the knowledge base*

905  A web based browser was created for browsing the knowledge base with Python, Flask, and
906  SQLAlchemy and hosted on Google Compute Engine, herein referred to Molecular Oncology
907  Almanac Browser or browser. The front page lists the total number of molecular features and
908  assertions catalogued as well as the total number of cancer types, evidence levels, and
909  therapies entered. A central search box allows for searching across multiple search terms such
910  as evidence, gene, feature types, or feature type attributes (protein changes, genomic positions,
911  etc.). The browser also features an about page, which contains a hyperlink to download the
912  contents of the knowledge base. Users may submit entries for consideration into the database
913  with a web form, accessible through the "Submit entry" menu item.

914  *Application Program Interface (API)*

915  To interact with the knowledge base programmatically, an application program interface (API)
916  was built using Python and Flask to interface with the browser's underlying data structure.
917  Several get requests are available to list therapies, evidence levels, or genes as well as the
918  ability to get all or by id assertions, sources, feature definitions, features, feature attribute
919  definitions, or feature attributes. A post request is available to suggest a new assertion to the
920  database.

921  *Reducing the burden of crowdsourcing*

922  To reduce the burden of crowdsourcing, we created a Google Chrome extension, herein
923  referred to as Molecular Oncology Almanac Connector or connector, with Python and Flask.
924  The connector allows users to submit a DOI along with a feature type, cancer type, evidence
925  level, and therapy if relevant. The user's email address is also requested in order to follow up
926  about the nominated assertion. This is accomplished using the post request API endpoint for
927  new assertions. The privacy policy of the Connector was reviewed and approved by Dana-
928  Farber compliance.

929  *Creating a cloud-based execution portal*

930  A web portal was built using Python, Flask, and requests to take advantage of Terra's (formerly
931  known as FireCloud) API and Google Cloud's gsutil in order to allow run MOAlmanac without
932  needing to use Python, Github, Docker, or Terra. Users must have billing set up with and be
933  registered on Terra and, upon selecting to begin a new analysis, users will be asked to specify a
934  de-identified sample name, either a free text tumor type or select one based on a drop down
935  menu containing ontologies from Oncotree, and a Terra billing project. A workspace will be
936  created in the specified billing project named based on the sample name, tumor type, and a
937  timestamp. The remaining fields are optional and any combination of them can be provided.

938 Somatic single nucleotide variants, insertions and deletions, bases covered, copy number
939 alterations, fusions, and somatic variants from orthogonal sequencing as well as a free text
940 description can be uploaded to the workspace through the web portal. The privacy policy and
941 application were reviewed and approved by Dana-Farber compliance and information security;
942 Nonetheless, we decided to remove germline inputs via the portal.
943
944 Upon submission, a Terra workspace and corresponding Google bucket is created that only the
945 user has access to and provided files are uploaded to the Google bucket. The workspace and
946 data model are populated based on inputs and a submission of Molecular Oncology Almanac is
947 run. The workspace is tagged with the tag Molecular-Oncology-Almanac-Portal on Terra. The
948 user is returned to their homepage on the portal, showing a summary of workspaces submitted
949 through the portal, by subsetting workspaces that they have access for the portal's tag. The
950 summary will note the job submission until the page. Upon page refresh with the job being
951 completed, a direct hyperlink to view the report output (View Report) is made available.

952 *Analysis and data availability*

953 All analyses and figures referenced herein can be found in and regenerated with the paper's
954 Github repository: https://github.com/brendanreardon/moalmanac-paper. Code is available for
955 all software in the Molecular Oncology Almanac ecosystem: browser
956 (https://github.com/vanallenlab/almanac-browser), connector (Google Chrome extension,
957 https://github.com/vanallenlab/almanac-extension), method
958 (https://github.com/vanallenlab/moalmanac), and portal
959 (https://github.com/vanallenlab/almanac-portal).
960
961
962
963

## Acknowledgements

971    Author information

972    **Affiliations**

973    **Department of Medical Oncology, Dana-Farber Cancer Institute, Harvard Medical School,**
974    **Boston, MA, USA**
975    Brendan Reardon, Nathaniel D Moore, Nicholas Moore, Eric Kofman, Saud Aldubayan,
976    Alexander Cheung, Jake Conway, Haitham Elmarakeby, Tanya Keenan, Daniel Keliher, David
977    Liu, Jihye Park, Natalie Vokes, Felix Dietlein, Eliezer M Van Allen
978
979    **Cancer Program, Broad Institute of MIT and Harvard, Cambridge, MA, USA**
980    Brendan Reardon, Nathaniel D Moore, Nicholas Moore, Eric Kofman, Saud Aldubayan,
981    Alexander Cheung, Jake Conway, Haitham Elmarakeby, Alma Imamovic, Sophia C. Kamran,
982    Tanya Keenan, Daniel Keliher, David J Konieczkowski, David Liu, Kent Mouw, Jihye Park,
983    Natalie Vokes, Felix Dietlein, Eliezer M Van Allen
984
985    **Indiana University School of Medicine, Indianapolis, IN, USA**
986    Nathaniel D Moore
987
988    **Howard Hughes Medical Institute, Chevy Chase, MD, USA**
989    Nathaniel D Moore
990
991    **Department of Internal Medicine, University of Cincinnati, Cincinnati, Ohio, USA**
992    Nathaniel D Moore
993
994    **Harvard Medical School, Harvard University, Boston, MA, USA**
995    Nicholas Moore, Kent Mouw
996
997    **Department of Cellular and Molecular Medicine, University of California, San Diego, La**
998    **Jolla, CA, USA**
999    Eric Kofman
1000
1001   **Institute for Genomic Medicine, University of California, San Diego, La Jolla, CA, USA**
1002   Eric Kofman
1003
1004   **Division of Genetics, Brigham and Women's Hospital, Boston, MA, USA**
1005   Saud Aldubayan
1006
1007   **College of Medicine, King Saud bin Abdulaziz University for Health Sciences, Riyadh,**
1008   **Saudi Arabia**
1009   Saud Aldubayan
1010
1011   **Grossman School of Medicine, New York University, New York, NY, USA**
1012   Alexander Cheung

1013

1014 **Division of Medical Sciences, Harvard University, Boston, MA, USA**

1015 Jake Conway

1016

1017 **Department of System and Computer Engineering, Al-Azhar University, Cairo, Egypt**

1018 Haitham Elmarakeby

1019

1020 **Department of Pediatric Oncology, Dana-Farber Cancer Institute, Harvard Medical**

1021 **School, Boston, MA, USA**

1022 Alma Imamovic

1023

1024 **Department of Radiation Oncology, Massachusetts General Hospital, Harvard Medical**

1025 **School, Boston, MA, USA**

1026 Sophia Kamran

1027

1028 **Department of Mathematics, Tufts University, Medford, MA, USA**

1029 Daniel Keliher

1030

1031 **Department of Radiation Oncology, Dana-Farber Cancer Institute & Brigham and**

1032 **Women's Hospital, Boston, MA**

1033 David J Konieczkowski, Kent Mouw

1034

1035 **Harvard Radiation Oncology Program, Massachusetts General Hospital, Boston, MA,**

1036 **USA**

1037 David J Konieczkowski

1038

1039 **Department of Radiation Oncology, The Ohio State University Comprehensive Cancer**

1040 **Center - Arthur G. James Cancer Hospital and Richard J Solove Research Institute,**

1041 **Columbus, OH, USA**

1042 David J Konieczkowski

1043

1044 **Department of Thoracic / Head and Neck Oncology, MD Anderson Cancer Center,**

1045 **Houston, TX, USA**

1046 Natalie Vokes

1047 **Contributions**

1048 Conception and designs: B.R., N.D.M., N.M., E.K., F.D., E.M.V.A. Development of methodology:
1049 B.R., N.D.M., N.M., E.K., S.A., A.C., J.C., H.E., A.I., S.C.K., T.K., D.K., D.J.K., D.L., K.W., J.P.,
1050 N.V., F.D., E.M.V.A. Analysis and interpretation of data: B.R., N.D.M., N.M., E.K., E.M.V.A.
1051 Writing, review, and/or revision of the manuscript: B.R., N.D.M., N.M., E.K., S.A., A.C., J.C.,
1052 H.E., A.I., S.C.K., T.K., D.K., D.J.K., D.L., K.W., J.P., N.V., F.D., E.M.V.A. Study supervision:
1053 E.M.V.A.

1054 **Competing interest statement**

1055 E.M.V.A. holds consulting roles with Tango Therapeutics, Genome Medical, Invitae, Enara Bio,

1056 Janssen, Manifold Bio, Monte Rosa. E.M.V.A. has received research support from Novartis,

1057 BMS. E.M.V.A. owns equity in Tango Therapeutics, Genome Medical, Syapse, Enara Bio,

1058 Manifold Bio, Microsoft, and Monte Rosa and has received travel reimbursement from

1059 Roche/Genentech. E.M.V.A., B.R., and N.D.M. have institutional patents filed on methods for

1060 clinical interpretation.


1061 **Corresponding author**

1062 Correspondence to Eliezer M Van Allen (EliezerM_VanAllen@dfci.harvard.edu)
1063

## References

1.  AACR Project GENIE Consortium. AACR Project GENIE: Powering Precision Medicine through an International Consortium. *Cancer Discov.* **7**, 818–831 (2017).

2.  Van Allen, E. M. *et al.* Whole-exome sequencing and clinical interpretation of formalin-fixed, paraffin-embedded tumor samples to guide precision cancer medicine. *Nat. Med.* **20**, 682–688 (2014).

3.  Chakravarty, D. *et al.* OncoKB: A Precision Oncology Knowledge Base. *JCO Precis Oncol* **2017**, (2017).

4.  Griffith, M. *et al.* CIViC is a community knowledgebase for expert crowdsourcing the clinical interpretation of variants in cancer. *Nat. Genet.* **49**, 170–174 (2017).

5.  Wagner, A. H. *et al.* A harmonized meta-knowledgebase of clinical interpretations of somatic genomic variants in cancer. *Nat. Genet.* **52**, 448–457 (2020).

6.  Patterson, S. E., Statz, C. M., Yin, T. & Mockus, S. M. Utility of the JAX Clinical Knowledgebase in capture and assessment of complex genomic cancer data. *NPJ Precis Oncol* **3**, 2 (2019).

7.  Tamborero, D. *et al.* Cancer Genome Interpreter annotates the biological and clinical relevance of tumor alterations. *Genome Med.* **10**, 25 (2018).

8.  Huang, K.-L. *et al.* Pathogenic Germline Variants in 10,389 Adult Cancers. *Cell* **173**, 355–370.e14 (2018).

9.  Polak, P. *et al.* A mutational signature reveals alterations underlying deficient homologous recombination repair in breast cancer. *Nat. Genet.* **49**, 1476–1486 (2017).

10. Larotrectinib OK'd for Cancers with TRK Fusions. *Cancer Discov.* **9**, 8–9 (2019).

11. Yizhak, K. *et al.* RNA sequence analysis reveals macroscopic somatic clonal expansion across normal tissues. *Science* **364**, (2019).

12. Van Hoeck, A., Tjoonk, N. H., van Boxtel, R. & Cuppen, E. Portrait of a cancer: mutational

1090        signature analyses for cancer diagnostics. *BMC Cancer* **19**, 457 (2019).

1091    13. Barretina, J. *et al.* 22 The Cancer Cell Line Encyclopedia - Using Preclinical Models to

1092        Predict Anticancer Drug Sensitivity. *European Journal of Cancer* vol. 48 S5–S6 (2012).

1093    14. Tsherniak, A. *et al.* Defining a Cancer Dependency Map. *Cell* **170**, 564–576.e16 (2017).

1094    15. Sinha, R., Schultz, N. & Sander, C. Comparing cancer cell lines and tumor samples by

1095        genomic profiles. *bioRxiv* 028159 (2015) doi:10.1101/028159.

1096    16. Najgebauer, H. *et al.* CELLector: Genomics-Guided Selection of Cancer In Vitro Models.

1097        *Cell Syst* **10**, 424–432.e6 (2020).

1098    17. Warren, A. *et al.* Global computational alignment of tumor and cell line transcriptional

1099        profiles. *bioRxiv* 2020.03.25.008342 (2020) doi:10.1101/2020.03.25.008342.

1100    18. Chang, M. T. *et al.* Accelerating Discovery of Functional Mutant Alleles in Cancer. *Cancer*

1101        *Discov.* **8**, 174–183 (2018).

1102    19. Babaei, S., Akhtar, W., de Jong, J., Reinders, M. & de Ridder, J. 3D hotspots of recurrent

1103        retroviral insertions reveal long-range interactions with cancer genes. *Nat. Commun.* **6**,

1104        6381 (2015).

1105    20. Gao, J. *et al.* 3D clusters of somatic mutations in cancer reveal numerous rare mutations as

1106        functional targets. *Genome Med.* **9**, 4 (2017).

1107    21. Sondka, Z. *et al.* The COSMIC Cancer Gene Census: describing genetic dysfunction

1108        across all human cancers. *Nat. Rev. Cancer* **18**, 696–705 (2018).

1109    22. Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set

1110        collection. *Cell Syst* **1**, 417–425 (2015).

1111    23. Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids*

1112        *Res.* **47**, D941–D947 (2019).

1113    24. Landrum, M. J. *et al.* ClinVar: improving access to variant interpretations and supporting

1114        evidence. *Nucleic Acids Res.* **46**, D1062–D1067 (2018).

1115    25. Karczewski, K. J. *et al.* The ExAC browser: displaying reference data information from over

1116      60 000 exomes. *Nucleic Acids Res.* **45**, D840–D845 (2017).

1117   26. Van Allen, E. M. *et al.* Genomic correlates of response to CTLA-4 blockade in metastatic

1118      melanoma. *Science* **350**, 207–211 (2015).

1119   27. Robinson, D. *et al.* Integrative Clinical Genomics of Advanced Prostate Cancer. *Cell* **162**,

1120      454 (2015).

1121   28. Bielski, C. M. *et al.* Genome doubling shapes the evolution and prognosis of advanced

1122      cancers. *Nat. Genet.* **50**, 1189–1195 (2018).

1123   29. Alexandrov, L. B., Nik-Zainal, S., Siu, H. C., Leung, S. Y. & Stratton, M. R. A mutational

1124      signature in gastric cancer suggests therapeutic strategies. *Nat. Commun.* **6**, 8683 (2015).

1125   30. Sztupinszki, Z. *et al.* Detection of molecular signatures of homologous recombination

1126      deficiency in prostate cancer with or without BRCA1/2 mutations. *Clin. Cancer Res.* (2020)

1127      doi:10.1158/1078-0432.CCR-19-2135.

1128   31. Chatterjee, P. *et al.* PARP inhibition sensitizes to low dose-rate radiation TMPRSS2-ERG

1129      fusion gene-expressing and PTEN-deficient prostate cancer cells. *PLoS One* **8**, e60408

1130      (2013).

1131   32. Ghandi, M. *et al.* Next-generation characterization of the Cancer Cell Line Encyclopedia.

1132      *Nature* **569**, 503–508 (2019).

1133   33. Yang, W. *et al.* Genomics of Drug Sensitivity in Cancer (GDSC): a resource for therapeutic

1134      biomarker discovery in cancer cells. *Nucleic Acids Res.* **41**, D955–61 (2013).

1135   34. Wang, B. *et al.* Similarity network fusion for aggregating data types on a genomic scale.

1136      *Nat. Methods* **11**, 333–337 (2014).

1137   35. Sicklick, J. K. *et al.* Molecular profiling of cancer patients enables personalized combination

1138      therapy: the I-PREDICT study. *Nat. Med.* **25**, 744–750 (2019).

1139   36. Lindsay, J. *et al.* MatchMiner: An open source computational platform for real-time

1140      matching of cancer patients to precision medicine clinical trials using genomic and clinical

1141      criteria. *bioRxiv* 199489 (2017) doi:10.1101/199489.

1142  37. Pallarz, S. *et al.* Comparative Analysis of Public Knowledge Bases for Precision Oncology.

1143  *JCO Precision Oncology* 1–8 (2019).

1144  38. Pai, S. & Bader, G. D. Patient Similarity Networks for Precision Medicine. *J. Mol. Biol.* **430**,

1145  2924–2938 (2018).

1146  39. Zitnik, M. *et al.* Machine Learning for Integrating Data in Biology and Medicine: Principles,

1147  Practice, and Opportunities. *Inf. Fusion* **50**, 71–91 (2019).

1148  40. Ramos, A. H. *et al.* Oncotator: cancer variant annotation tool. *Hum. Mutat.* **36**, E2423–9

1149  (2015).

1150  41. Lichtenstein, L., Woolf, B., MacBeth, A., Birsoy, O. & Lennon, N. Abstract 3641:

1151  ReCapSeg: Validation of somatic copy number alterations for CLIA whole exome

1152  sequencing. *Cancer Res.* **76**, 3641–3641 (2016).

1153  42. Haas, B. J. *et al.* Accuracy assessment of fusion transcript detection via read-mapping and

1154  de novo fusion transcript assembly-based methods. *Genome Biol.* **20**, 213 (2019).

1155  43. Kurian, A. W. *et al.* Clinical evaluation of a multiple-gene sequencing panel for hereditary

1156  cancer risk assessment. *J. Clin. Oncol.* **32**, 2001–2009 (2014).

1157  44. Kalia, S. S. *et al.* Recommendations for reporting of secondary findings in clinical exome

1158  and genome sequencing, 2016 update (ACMG SF v2.0): a policy statement of the

1159  American College of Medical Genetics and Genomics. *Genet. Med.* **19**, 249–255 (2017).

1160  45. Chen, E. J. *et al.* Abiraterone treatment in castration-resistant prostate cancer selects for

1161  progesterone responsive mutant androgen receptors. *Clin. Cancer Res.* **21**, 1273–1280

1162  (2015).

1163  46. Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-

1164  associated genes. *Nature* **499**, 214–218 (2013).

1165  47. Rosenthal, R., McGranahan, N., Herrero, J., Taylor, B. S. & Swanton, C. DeconstructSigs:

1166  delineating mutational processes in single tumors distinguishes DNA repair deficiencies

1167  and patterns of carcinoma evolution. *Genome Biol.* **17**, 31 (2016).

1168    48.  Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J. & Stratton, M. R.

1169         Deciphering signatures of mutational processes operative in human cancer. *Cell Rep.* **3**,

1170         246–259 (2013).

1171    49.  Salipante, S. J., Scroggins, S. M., Hampel, H. L., Turner, E. H. & Pritchard, C. C.

1172         Microsatellite instability detection by next generation sequencing. *Clin. Chem.* **60**, 1192–

1173         1199 (2014).

1174    50.  Maruvka, Y. E. *et al.* Analysis of somatic microsatellite indels identifies driver events in

1175         human tumors. *Nat. Biotechnol.* **35**, 951–959 (2017).

1176    51.  Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and

1177         heterogeneous cancer samples. *Nat. Biotechnol.* **31**, 213–219 (2013).

1178    52.  Saunders, C. T. *et al.* Strelka: accurate somatic small-variant calling from sequenced

1179         tumor–normal sample pairs. *Bioinformatics* **28**, 1811–1817 (2012).

1180    53.  Costello, M. *et al.* Discovery and characterization of artifactual mutations in deep coverage

1181         targeted capture sequencing data due to oxidative DNA damage during sample

1182         preparation. *Nucleic Acids Res.* **41**, e67 (2013).

1183    54.  Lawrence, M. S. *et al.* Discovery and saturation analysis of cancer genes across 21 tumour

1184         types. *Nature* **505**, 495–501 (2014).

1185    55.  Poplin, R. *et al.* A universal SNP and small-indel variant caller using deep neural networks.

1186         *Nat. Biotechnol.* **36**, 983–987 (2018).

1187    56.  Olshen, A. B., Venkatraman, E. S., Lucito, R. & Wigler, M. Circular binary segmentation for

1188         the analysis of array-based DNA copy number data. *Biostatistics* **5**, 557–572 (2004).

1189    57.  Shen, R., Olshen, A. B. & Ladanyi, M. Integrative clustering of multiple genomic data types

1190         using a joint latent variable model with application to breast and lung cancer subtype

1191         analysis. *Bioinformatics* **25**, 2906–2912 (2009).

1192    58.  Shen, R. & Seshan, V. E. FACETS: allele-specific copy number and clonal heterogeneity

1193         analysis tool for high-throughput DNA sequencing. *Nucleic Acids Res.* **44**, e131 (2016).

1194    59. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).

1195    60. Hass, B. *et al.* STAR-fusion: fast and accurate fusion transcript detection from RNA-Seq.

1196        bioRxiv. (2017).

1197    61. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing

1198        next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).

1199    62. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-

1200        generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).

1201    63. Van der Auwera, G. A. *et al.* From FastQ data to high confidence variant calls: the Genome

1202        Analysis Toolkit best practices pipeline. *Curr. Protoc. Bioinformatics* **43**, 11.10.1–11.10.33

1203        (2013).

1204    64. de Weck, A., Bitter, H. & Kauffmann, A. Fibroblasts cell lines misclassified as cancer cell

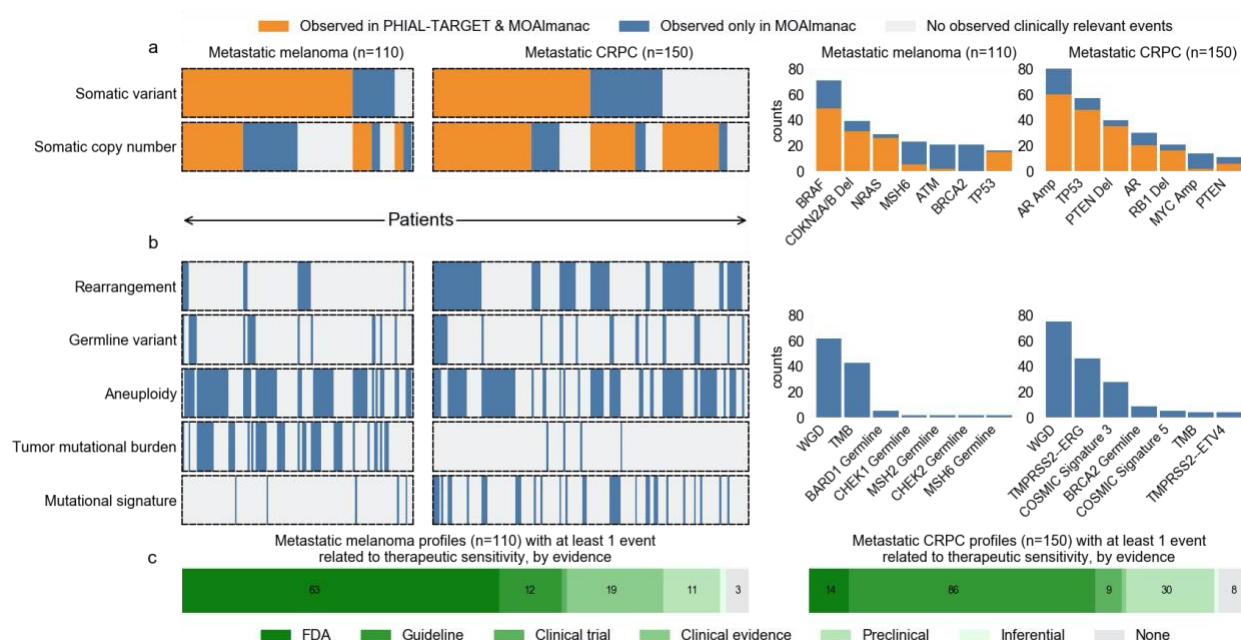1205        lines. *bioRxiv* 166199 (2017) doi:10.1101/166199.
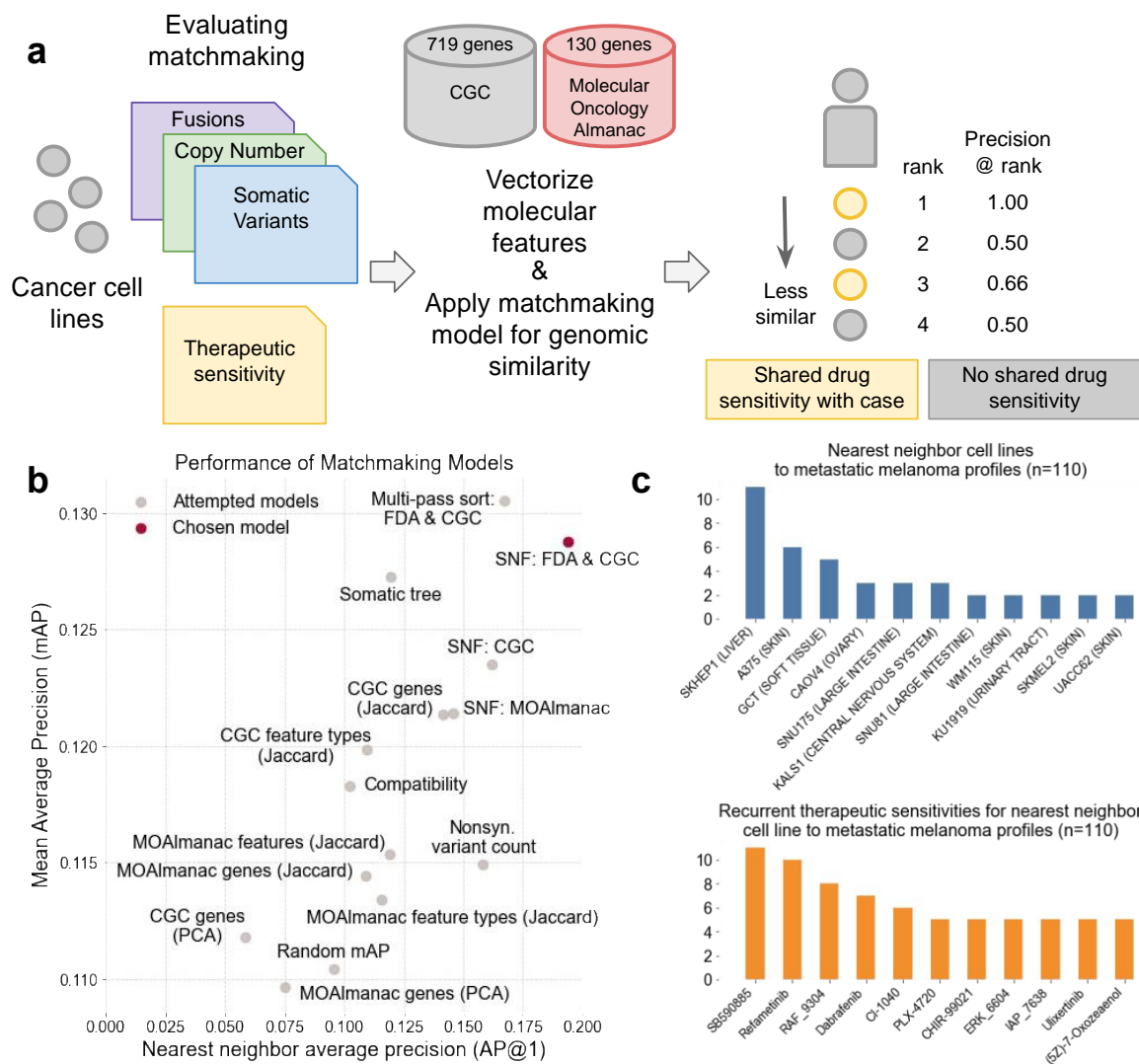
1206  Figures and captions



1207

1208  **Figure 1.** Molecular Oncology Almanac, a clinical interpretation framework

1209  **(a)** The Molecular Oncology Almanac accepts any combination of somatic single nucleotide
1210  variants (snvs), insertions and deletions (indels), copy number alterations (cnas), germline snvs
1211  and indels, somatic snvs from orthogonal sequencing, and rearrangements from RNA.
1212  Molecular features are annotated for clinical relevance and with several other data sources
1213  before being heuristically sorted (first-order). Variants are used to evaluate genomic features;
1214  somatic-germline overlap, concordance of somatic variants with orthogonal sequencing,
1215  COSMIC mutational signature contributions, mutational burden, and MSI related variants
1216  (second-order). Somatic mutations, copy number alterations, and fusions are used to assess
1217  similarity to individual cell lines for further therapeutic sensitivity suggestions. A report of
1218  putative actionability is generated (Methods). **(b)** A literature review was performed to identify
1219  relationships between molecular alterations and clinical actions for precision oncology,
1220  beginning with relationships suggested in TARGET[2]. 63 genes were removed from TARGET
1221  due to insufficient evidence and 58 were retained. Clinical relationships were cataloged as
1222  suggesting therapeutic sensitivity, resistance, or prognostic value in an SQL database
1223  (Methods**)** and made available online (https://moalmanac.org). (**c**) Sources catalogued in the
1224  Molecular Oncology Almanac, categorized by evidence (left) and therapy types (right)

1225

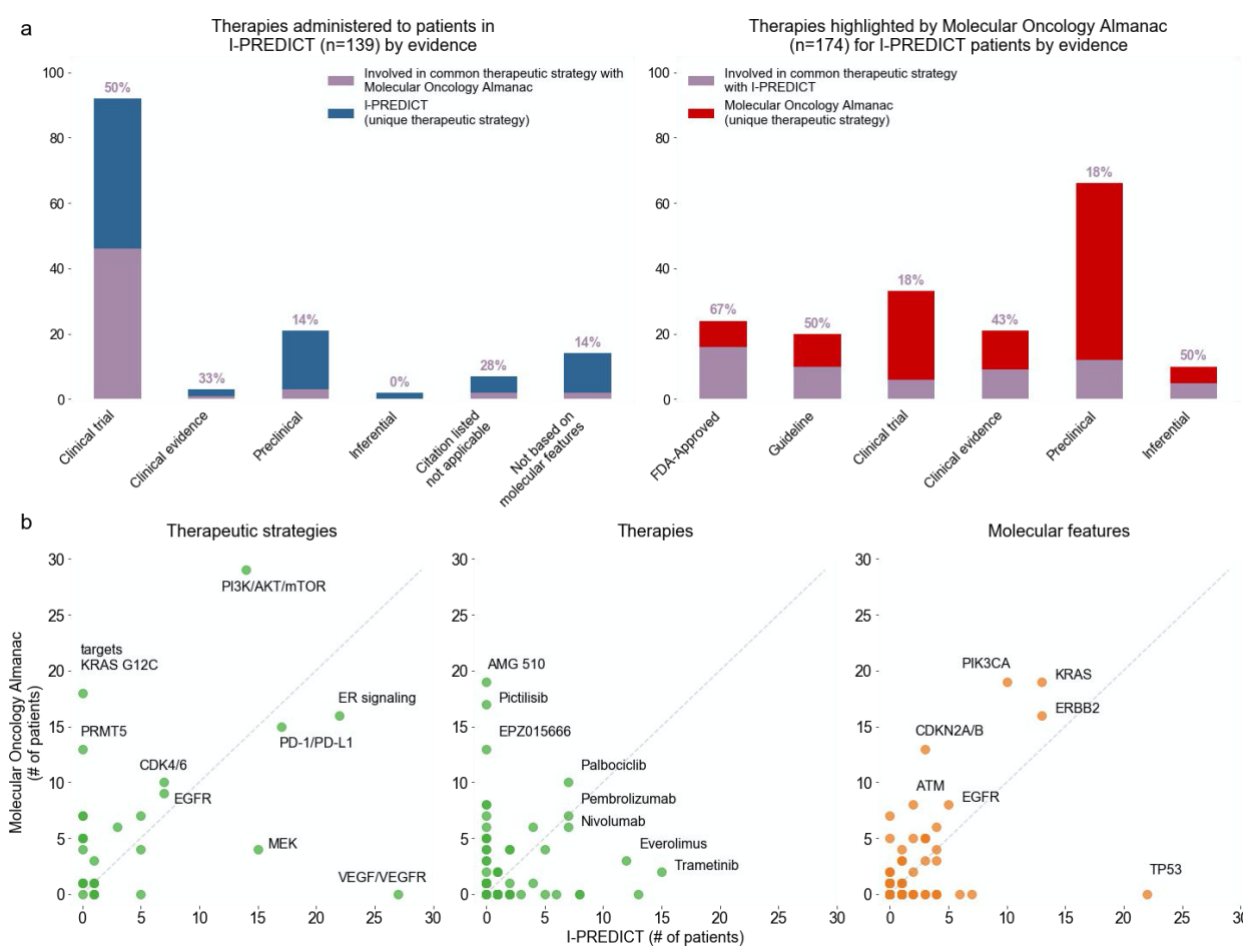**Figure 2.** Benchmarking MOAlmanac against PHIAL & TARGET.

Molecular Oncology Almanac was benchmarked against PHIAL & TARGET using 110 metastatic melanomas and 150 metastatic castration-resistant prostate cancers[2,26,27]. **(a)** The Molecular Oncology Almanac increased the number of patients with a somatic variant or copy number alteration labeled as "putatively actionable" or "investigate actionability" from 115 to 249 relative to PHIAL; patients are aligned across feature types vertically (left). Specific molecular features that were observed by both PHIAL & MOAlmanac (orange) and by MOAlmanac only (blue) for each cohort are shown (right). **(b)** Features not routinely used in clinical sequencing were utilized to characterize actionability: rearrangements, germline variants, aneuploidy, mutational burden, and mutational signatures; patients aligned with (A) vertically (left). Considering these features types further identified 7 patients with a clinically relevant feature. Specific molecular features that were additionally observed in each cohort are shown (right). (Abbreviations used: WGD = whole genome doubling, TMB = tumor mutational burden). (**c**) Including preclinical evidence when considering putative actionability provides an additional 41 patients (11 patients with metastatic melanoma and 30 patients with castration resistant prostate cancers) with a molecularly matched therapeutic hypothesis.

1242

**Figure 3.** Leveraging preclinical models in MOAlmanac.

MOAlmanac leverages preclinical data from cancer cell lines which have been molecularly characterized and subject to high-throughput therapeutic screens to provide supplemental hypotheses through profile-cell line matchmaking. (**a**) Somatic SNVs, CNAs, and fusions of cancer cell lines are formatted, annotated with MOAlmanac and CGC, and vectorized into sample x feature boolean dataframes. Feature sets and similarity metrics were evaluated by their ability to sort cell lines relative to one another based on shared genomic features, such that cell lines that shared therapeutic sensitivity were deemed more similar. Metrics from information retrieval were used for evaluation; mean average precision (mAP, how the model does overall at sorting cell lines which share therapeutic sensitivity to be closer to the case profile) and average precision at rank 1 (ap@1, how often the nearest neighbor shared therapeutic sensitivity). (**b**) Models were evaluated on 377 cancer cell lines using a hold-one-out approach. The model which had the strongest trade off between the two metrics used Similarity Network Fusion to fuse networks of somatic variants, copy number alterations, and fusions in CGC genes with specific MOAlmanac features associated with an FDA approval[21,34]. (**c**) Recurrent

1258    nearest neighbors and their sensitive therapies for 110 metastatic melanomas. SKHEP1_LIVER was

1259    the first neighbor for 11 profiles, A375_SKIN for six, and GCT_SOFT_TISSUE for five. Nearest

1260    neighbors were sensitive to MEK and RAF inhibitors: SB590885 (BRAF inhibitor, 11 neighbors),

1261    Refametinib (MEK, 10), RAF_9304 (RAF, 8), and Dabrafenib (BRAF, 7).

1262

**Figure 4.** Application of MOAlmanac to a prospective clinical trial.

We investigated if MOAlmanac could highlight similar therapeutic strategies that were utilized by real world evidence. MOAlmanac was applied to the I-PREDICT trial, which evaluated the efficacy of molecularly matched therapies in 83 patients[35]. (**a**) Therapies and corresponding molecular features were mapped to therapeutic strategies for those administered in I-PREDICT and highlighted by MOAlmanac. MOAlmanac nominated therapeutic strategies applied for a given patient (purple) more often for those based on well established evidence (i.e. FDA approvals; 67% of therapy patient pairs) relative to less established evidence, such as preclinical (18%). Counts of therapeutic strategies applied to patients that were unique to I-PREDICT are shown in blue and those highlighted by and unique to MOAlmanac are in red. (**b**) Therapeutic strategies, individual therapies, and molecular features as administered or targeted by I-PREDICT and highlighted by Molecular Oncology Almanac.