

1 **The gene expression classifier ALLCatchR identifies B-precursor ALL**
2 **subtypes and underlying developmental trajectories across age**

3 Thomas Beder¹, Björn-Thore Hansen¹, Alina M. Hartmann^{1,2}, Johannes Zimmermann³, Eric
4 Amelunxen¹, Nadine Wolgast^{1,2}, Wencke Walter⁴, Marketa Zaliova⁵, Željko Antić⁶, Philippe
5 Chouvarine⁶, Lorenz Bartsch¹, Malwine Barz^{1,2}, Miriam Bultmann¹, Johanna Horns¹, Sonja
6 Bendig^{1,2}, Jan Kässens¹, Christoph Kaleta³, Gunnar Cario^{2,7}, Martin Schrappe^{2,7}, Martin
7 Neumann^{1,2}, Nicola Gökbüget⁸, Anke Katharina Bergmann⁶, Jan Trka⁵, Claudia Haferlach⁴,
8 Monika Brüggemann^{1,2}, Claudia D. Baldus^{1,2}, Lorenz Bastian^{1,2}

9

10 1 Medical Department II, Hematology and Oncology, University Hospital Schleswig-Holstein,
11 Kiel, Germany

12 2 Clinical Research Unit 'CATCH ALL' (KFO 5010/1) funded by the Deutsche
13 Forschungsgemeinschaft

14 3 Institute of Experimental Medicine, Research Group Medical Systems Biology, Christian-
15 Albrechts-University Kiel, Germany

16 4 MLL Munich Leukemia Laboratory, Munich, Germany

17 5 Childhood Leukaemia Investigation Prague, and University Hospital Motol, Prague, Czech
18 Republic

19 6 Department of Human Genetics, Hannover Medical School (MHH), Hannover, Germany

20 7 Department of Pediatrics, University Hospital Schleswig-Holstein Kiel, Kiel, Germany

21 8 Department of Medicine II, Hematology/Oncology, Goethe University Hospital,
22 Frankfurt/M., Germany

23

24 Short title: **ALLCatchR identifies B cell precursor-ALL subtypes**

25 Corresponding author: Claudia Baldus

26 Medical Department II, Hematology and Oncology

27 Universitätsklinikum Schleswig-Holstein, Campus Kiel

28 Arnold-Heller-Strasse 3

29 24105 Kiel, Germany

30 Mail: claudia.baldus@uksh.de

31 Phone: +49 431 500 22500

32

33 Word count: 4000

34 References: 35

35 Figures: 5

36

37 **Abstract**

38 Current classifications (WHO-HAEM5 / ICC) define up to 26 molecular B-cell precursor acute
39 lymphoblastic leukemia (BCP-ALL) disease subtypes, which are defined by genomic driver
40 aberrations and corresponding gene expression signatures. Identification of driver aberrations
41 by RNA-Seq is well established, while systematic approaches for gene expression analysis are
42 less advanced. Therefore, we developed ALLCatchR, a machine learning based classifier using
43 RNA-Seq expression data to allocate BCP-ALL samples to 21 defined molecular subtypes.
44 Trained on n=1,869 transcriptome profiles with established subtype definitions (4 cohorts;
45 55% pediatric / 45% adult), ALLCatchR allowed subtype allocation in 3 independent hold-out
46 cohorts (n=1,018; 75% pediatric / 25% adult) with 95.7% accuracy (averaged sensitivity across
47 subtypes: 91.1% / specificity: 99.8%). ‘High confidence predictions’ were achieved in 84.6% of
48 samples with 99.7% accuracy. Only 1.2% of samples remained ‘unclassified’. ALLCatchR
49 outperformed existing tools and identified novel candidates in previously unassigned samples.
50 We established a novel RNA-Seq reference of human B-lymphopoiesis. Implementation in
51 ALLCatchR enabled projection of BCP-ALL samples to this trajectory, which identified shared
52 patterns of proximity of BCP-ALL subtypes to normal lymphopoiesis stages. ALLCatchR sustains
53 RNA-Seq routine application in BCP-ALL diagnostics with systematic gene expression analysis
54 for accurate subtype allocations and novel insights into underlying developmental
55 trajectories.

56

57 Introduction

58 Improved outcomes in B cell precursor acute lymphoblastic leukemia (BCP-ALL) – both, in
59 pediatric and adult patients – have been achieved by precise risk stratification and target
60 specific treatments. Molecular BCP-ALL subtypes and immunophenotype are the most
61 important baseline prognosticators for BCP-ALL beside white blood cell counts and age. They
62 inform risk-adapted treatments and targeted therapies. Currently, the revised WHO
63 classification of lymphoid neoplasms, (WHO-HAEM5)¹ and the International Consensus
64 Classification of Myeloid Neoplasms and Acute Leukemia (ICC)² have acknowledged 11 and 26
65 molecular defined BCP-ALL subtypes as distinct diagnostic entities, respectively, including 5
66 provisional entities (ICC classification). A total of 21 of these subtypes have been characterized
67 by distinct gene expression profiles³⁻⁸, while the remaining subtypes^{2,5} are rare (IGH::*IL3*) or
68 were defined by specific sets of underlying genomic drivers (Ph-like: ABL class / JAK-STAT /
69 NOS) or their absence (*KMT2A-/ZNF384*-like). This heterogeneity of diagnostic subtypes
70 exceeds the capabilities of cytogenetic (chromosome banding analysis, FISH) and molecular
71 genetic methods (breakpoint specific PCR, MLPA, SNP-array / Array-CGH) which have so far
72 been combined for identification of BCP-ALL subtypes. RNA-Seq enables identification of all
73 BCP-ALL subtypes with a single method, establishing a new diagnostic standard. Further
74 implementation as routine clinical diagnostic requires unified analysis methods. Calling of
75 driver gene fusions^{9,10} is well established and novel approaches for the identification of
76 hotspot single nucleotide¹⁰ variants and virtual karyotypes¹¹ exist. Yet only few approaches
77 for systematic gene expression analysis are currently available.

78 Gene expression signatures represent the signaling equivalent of heterogeneous genomic
79 driver alterations and have been used to define BCP-ALL subtypes. Initially, unsupervised
80 clustering or prediction analysis for microarrays (PAM) were used to define subtype specific
81 gene sets resulting in considerable heterogeneity regarding gene set definitions and subtype
82 allocation of individual samples.¹² More recent systematic approaches for BCP-ALL subtype
83 allocations have employed machine learning methods to train classifiers for BCP-ALL subtype
84 allocation mainly on pediatric ALL datasets.^{13,14} Yet the optimal method still needs to be
85 defined – especially for rare and difficult to classify subtypes and subtypes with predominance
86 in adults. Additionally, correct assignment of samples, which do not fall into established
87 subtype categories either due to interfering biological conditions (e.g., low blast count, poor
88 RNA quality) or because these samples represent novel candidate subtypes, remains a

89 challenge. In addition to molecular subtype definitions, gene expression profiles might be
90 informative for clinical baseline parameters such as leukemic blast proportion,
91 immunophenotype or more detailed analysis of lymphopoiesis trajectories underlying BCP-
92 ALL development. However, systematic approaches and especially RNA-Seq data that link BCP-
93 ALL subtypes to human B lymphopoiesis differentiation stages are lacking.

94 Here we describe ALLCatchR, a machine learning based classifier pretrained for allocation of
95 BCP-ALL gene expression profiles to all 21 gene expression defined molecular subtypes of
96 WHO-HAEM5 and ICC classifications. High accuracies in independent validation cohorts are
97 achieved by integrating machine learning and gene set based nearest neighbor models into a
98 compound classifier. ALLCatchR infers clinical baseline variables such as blast proportion and
99 patient's sex from RNA-Seq data and provides a putative differentiation stage of origin based
100 on our newly established reference of human B lymphopoiesis. ALLCatchR sustains routine
101 diagnostic application of RNA-Seq with systematic gene expression analysis providing subtype
102 allocations and insights into underlying biology for further exploratory analysis.

103 **Material/Subjects and Methods**

104 **Aggregation of a 3,532 sample BCP-ALL transcriptome reference data set.**

105 To establish a classifier for BCP-ALL molecular subtype allocation, we aggregated RNA-Seq
106 count data from n= 3,532 BCP-ALL patients including 64.5% pediatric^{5-7,13} and 35.5% adult<sup>3-
107 5,8,13</sup> cases combined from 6 independent datasets (**Figure 1A; Supplementary Table S1**).
108 Excluded were samples with multiple subtype assignments (n=116), multiple representations
109 of the same patient (n=44), subtypes which are not part of WHO-HAEM5 / ICC classification
110 (Low hyperdiploid, *IDH1/2*; n=55) or which are mainly defined by absence of a genomic driver
111 (*KMT2A*-like, *ZNF384*-like; n=9). Molecular BCP-ALL subtype allocations were performed for
112 n=2,887 samples in the original studies based on genomic drivers and corresponding gene
113 expression signatures. Subtype-defining genomic events were identified in >90% of cases
114 either by RNA-Seq (gene fusions, hotspot single nucleotide variants, virtual karyotypes) or by
115 genomic profiling (whole genome- / whole exome- / gene panel sequencing, SNP-arrays,
116 array-CGH). A total of n=421 samples were defined 'unassigned' or 'B-other' in the original
117 studies. All BCP-ALL molecular subtypes from current WHO-HAEM5 or ICC classifications which
118 were characterized by distinct gene expression signatures in their original description (n=21)

119 were represented in the data set (not included: *IGH::IL3*, *KMT2A*-like, *ZNF384*-like. Ph-like was
120 considered one subtype without sub-division. *CEBP/ZEB2* subtype lacks final definitions so far
121 and was defined here as 'CEBP' by presence of *IGH::CEBPA/CEBPE/CEBPD* fusions and absence
122 of other drivers.) (**Supplementary Table S2**). Raw read counts for 15,728 protein-coding genes
123 represented in all cohorts were used including heterogenous sequencing approaches (poly-A
124 selection / depletion of ribosomal RNAs), sequencing depths and different read count
125 quantification methods before normalization ($\log_{10}(\text{count} + 1)$), followed by z-transformation
126 and scaling between 0-1). The data set was split into a data set used for training of the classifier
127 ($n=1,869$) and 3 hold-out studies ($n=1,018$) for independent validation both representing all
128 analyzed BCP-ALL subtypes (**Figure 1A, Supplementary Figure S1**).

129 **Integration of machine learning and gene set based nearest-neighbor models for BCP-ALL** 130 **subtype allocation**

131 To perform molecular subtype allocation based exclusively on gene expression data, we
132 developed ALLCatchR, a classifier which integrates linear support vector machine (SMV) and
133 nearest-neighbor association models for BCP-ALL subtypes derived from the training data
134 (**Supplementary Figure S1**) Feature selection (LASSO)¹⁵ using the glmnet package¹⁶ was used
135 to extract BCP-ALL subtype defining gene sets, resulting in 2,802 genes with high
136 discriminative power for 21 molecular subtypes (**Supplementary Figure S2, Supplementary**
137 **Table S3**). First, we used this gene set to train five different machine learning classifiers using
138 two feature selection methods^{15,17} of which linear SVM¹⁸ performed best independent of the
139 feature selection method used (**Supplementary Figure S3**). This resulted in a high accuracy
140 (0.963) of subtype prediction in the training data. However, linear SVM is restricted to
141 predefined classes and does not compute probabilities for individual subtype predictions,
142 which prevents it from correctly handling cases which are unassigned or ambiguous due to
143 multiple drivers or which represent novel candidates. To achieve a probabilistic compound
144 model, we incorporated single sample gene set enrichment analyses (ssGSEA) using
145 singscore¹⁹ of the same subtype-defining LASSO gene sets. By this approach, batch effects
146 between cohorts were removed (**Supplementary Figure S4**) Euclidean distance of each test
147 sample to each training sample was computed and the 10 nearest neighbors were considered
148 for subtype allocations of each test sample (accuracy for subtype prediction based on highest
149 enrichment for each sample: 0.912). Both models - SVM linear predictions and sample-to-
150 samples-distances in subtype-defining gene sets – were integrated into our newly established

151 compound classifier, ALLCatchR, which provides dynamic ranges of subtype-specific
152 probability scores (**Figure 1A**). To achieve a better separation between highly similar high
153 hyperdiploid and near haploid ALL, both subtypes were first represented as one class in the
154 overall classifier (NH/HeH) and then separated by a second 2-class compound classifier with a
155 similar design as the overall classifier.

156 **Development of an RNA-Seq reference of human B-lymphopoiesis**

157 Bone marrow samples from healthy adult donors (n=4, M:F=1:3, age: 27-39 years, study
158 registration DRKS00023583, ethical approval of ethics committee, Kiel University: D 583/20)
159 were subjected to immunodensity cell separation (RosetteSep, STEMCELL Technologies; Inc.,
160 Vancouver, BC, Canada; purging: CD16, CD36, CD66b, CD235a, CD3). Non-depleted cells were
161 stained with a 9-color antibody panel (**Supplementary Table S4**) and FACS-sorted (FACSAria™
162 fusion; BD Biosciences, Franklin Lakes, NJ, USA) to 7 lymphoid differentiation stages. RNA was
163 extracted from 5,000-320,000 cells per differentiation stage (AllPrep™ DNA/RNA Micro Kit,
164 Qiagen, Venlo, Netherlands) and subjected to ultra-low-input RNA sequencing after
165 generation of stranded sequencing libraries (SMART-Seq® Stranded Kit, Takara Bio Inc.,
166 Kusatsu, Shiga, Japan; NovaSeq 6000, Illumina, San Diego, CA, USA).

167 **Results**

168 **ALLCatchR performs BCP-ALL molecular subtype allocation with high accuracy**

169 We used aggregated BCP-ALL gene expression profiles (n=3,532 samples, n=6 cohorts) to
170 develop ALLCatchR, a pre-trained machine learning classifier which performs BCP-ALL
171 molecular subtype allocation based on gene expression alone (detailed in 'Methods').
172 ALLCatchR provides probability scores for each sample and all gene expression defined BCP-
173 ALL subtypes (**Figure 1A**). Unsupervised clustering of ALLCatchR scores groups samples
174 according to subtype across cohorts and age groups. For final subtype allocation, we defined
175 subtype-specific cutoffs based on the comparison of probability scores from samples
176 belonging to the corresponding subtype and all remaining samples of the cohort (**Figure 1B**).
177 This resulted in 1.) high-confidence predictions, 2.) candidate predictions and 3.) low-
178 confidence predictions i.e., unclassified samples. Cutoffs for 'high-confidence' predictions
179 were defined to include >90% of correct predictions. Cutoffs for candidate predictions were

180 defined to exclude all samples from other subtypes but allowed unassigned/B-other samples
181 (n=111; **Figure 1B**). In the training data, 84.6% of samples achieved high confidence
182 predictions with an accuracy of 0.997, while 13.7% achieved candidate predictions with an
183 accuracy of 0.797 to guide further validation based on genomic drivers in well pre-specified
184 directions (**Figure 1C**). Only 1.7% of samples achieved low-confidence predictions and were
185 considered ‘unclassified’. To validate ALLCatchR performance, we used independent
186 validation data from 3 hold-out cohorts (n=1018; **Supplementary Figure S1A**), not previously
187 seen by the classifier. A total of n=1006 (98.8%) samples was allocated to one of 21 subtypes
188 (high-confidence and candidate predictions) with an accuracy of 0.957, demonstrating the
189 feasibility of highly accurate subtype allocations based on gene expression alone. ‘High-
190 confidence’ and ‘candidate’ predictions were achieved in 83.7% and 15.1% of samples with
191 accuracies of 0.989 and 0.851 respectively. A total of n=32 samples (3.1%) were assigned to
192 the wrong subtype or received no subtype allocation (n=12; 1.2%). Most prominent
193 misclassifications affected Ph-like- to Ph-pos predictions or vice versa (n=8) or subtype
194 allocations of aneuploid subtypes (n=20), **Figure 1D**). The majority of misclassified samples
195 (n=23/33; 67.6%) had received candidate predictions, supporting the need to validate these
196 predictions based on genomic drivers.

197 **ALLCatchR provides subtype allocations for previously ‘unassigned / B-other’ samples**

198 In addition to the n=1018 hold-out samples with assigned subtype, n=111 samples had been
199 defined as ‘unassigned / B-other’ (n=107) or were identified as ‘non-Ph-like *CRLF2*-rearranged’
200 (n=4) in the original studies (**Figure 1C, D**). ALLCatchR concordantly identified n=20 (18.0%) of
201 these as ‘unclassified’ (**Figure 1D, Supplementary Figure S5**). However, n=43 (38.7 %) and
202 n=48 (43.2 %) cases received ‘high-confidence’ or ‘candidate’ predictions respectively (**Figure**
203 **1D**). Analysis of available RNA-Seq gene fusion calls or cytogenetic profiles and/or virtual
204 karyotyping (WGS / SNP-arrays) identified driver candidates supporting the corresponding
205 subtype allocations in n=31 (72.1%) of ‘high-confidence’ and n=13 (27.1%) of ‘candidate’
206 predictions (**Supplementary Table S5; Supplementary Figure S5**). These newly suggested
207 subtype allocations consisted of PAX5alt predictions (n=25) which had not shown a clear
208 PAX5alt gene expression profile in the original cohort (n=1), or which were contributed from
209 the CLIP cohort where this subtype had not been annotated previously. Next, n=11 *CRLF2*-
210 rearranged cases from CLIP and St Jude cohorts without Ph-like gene expression profile in the
211 original cohorts received ALLCatchR Ph-like predictions. Among the remaining n=7 samples,

212 one case with an ALLCatchR high-confidence KMT2A prediction was found to harbor a *KMT2A*
213 partial tandem duplication by WGS (**Supplementary Figure S5**). To the best of our knowledge,
214 this is the first identification in BCP-ALL of this aberration which is recurrently observed in
215 acute myeloid leukemia. In a second of these n=7 cases, an *IGH::MYC* gene fusion was
216 identified in support of a *BCL2/MYC* ALLCatchR prediction. Further ALLCatchR high-confidence
217 predictions for ‘unassigned / B-other samples’ without corresponding drivers included PAX5alt
218 (n=9) and Ph-like (n=3) predictions, which generally are defined in a proportion of samples by
219 gene expression alone. Thus, ALLCatchR suggested molecular subtype allocations in previously
220 ‘unassigned’ cases with atypical and less well-defined gene expression signatures and
221 supported the identification of novel driver candidates.

222 **High accuracy of ALLCatchR predictions is observed across cohorts and molecular subtypes**

223 The accuracy of predictions was consistently high in the training and hold-out data, with 0.952
224 and 0.957, respectively. Almost congruent predictions were achieved in St Jude and CLIP
225 cohorts with accuracies of 0.978 and 0.965, respectively. In the MLL hold out set the accuracy
226 was slightly lower with 0.914 (**Figure 2A**). Of note, the MLL cohort includes real-world adult
227 BCP-ALL data from a diagnostic laboratory with more permissive pre-selection cutoffs (e.g.,
228 blast counts) indicating that ALLCatchR achieves reliable predictions also in less pre-selected
229 samples. Despite the overall high accuracies, classification performance varied between
230 molecular subtypes (**Figure 2B**). ALLCatchR achieved specificities >0.99 for all 21 subtypes,
231 both in training and testing data sets. The average sensitivity across subtypes was 0.919 ± 0.145
232 and 0.911 ± 0.167 in the training and hold-out data, respectively. For n=17/21 subtypes,
233 sensitivities were ≥ 0.85 both on training and hold-out data, together including n=2,781
234 patients (96.3%; **Figure 2B**). Only 4 remaining subtypes (n=106 samples, 3.7% of entire cohort)
235 achieved sensitivities below 0.85 (*NUTM1*, *CEBP*, *iAMP21* and Near haploid) which was mainly
236 related to the small number of samples representing these subtypes, limiting both.

237 **ALLCatchR subtype allocation outperforms current tools**

238 Recently, two tools - ALLSorts¹³ and Allspice¹⁴ - were independently developed for BCP-ALL
239 subtype allocation based on gene expression profiles. In comparison to these, ALLCatchR
240 provides comprehensive subtype-allocation to all gene expression defined WHO / ICC
241 subtypes (n=21), including *CEBP* and *CDX2/UBTF*, which are missed by both tools. For
242 performance comparison, n=2,887 samples with established subtype definitions (ALLCatchR

243 training and validation data sets) were predicted with ALLSorts and Allspice (**Supplementary**
244 **Figure S6**). ALLSorts performed well with an accuracy of 0.913 but left more samples
245 ‘unclassified’ (n=145), compared to ALLCatchR (n=44). The largest difference was observed in
246 the MLL holdout data, where ALLSorts achieved an accuracy of 0.771 compared to 0.914
247 accuracy for ALLCatchR (**Supplementary Figure S6**). An inferior performance in ALLSorts was
248 mainly related to missed sample classifications (ALLSorts ‘unassigned’: n=43 (16.17%);
249 ALLCatchR ‘unassigned’: n=8 (3.01%)). The MLL data set represents real-world data from a
250 diagnostic laboratory with less stringent pre-selection of samples and thus represent a *bona*
251 *fide* challenge for the tools. Allspice leaves more samples of the same cohort (n=2,887)
252 unclassified resulting in accuracies of 0.629 in the training and 0.719 in the hold-out studies.
253 However, for samples that could be assigned to a subtype by Allspice, prediction
254 performances were comparable to ALLCatchR (**Supplementary Figure S6**). In summary,
255 ALLCatchR achieves a higher accuracy for molecular subtype predictions, assigning more
256 samples to the correct subtype including all gene expression defined subtypes.

257 **Gene expression-based modelling predicts clinical baseline variables.**

258 Blast count proportions impact accuracy of gene expression based molecular subtype
259 allocation, as sequencing reads from non-leukemic compartments contribute to bulk
260 transcriptome profiles. To infer sample blast proportions, we trained two machine learning
261 regression models on data sets of our combined cohort with available blast counts obtained
262 by manual counting or flow cytometry (GMALL, MLL) and used these as well as the RCH/PM
263 cohort for validation. Blast count predictions from single cohorts achieved good accuracies
264 when applied to each other (**Figure 3A-B**) with a high concordance between USKH and MLL
265 training sets (**Figure 3B**) which were therefore combined for the final classifier. Only 1.85% of
266 samples with high confidence subtype predictions had blast count predictions <50% while
267 these were observed in 9.83% of candidate predictions and in 17.95% of unclassified samples
268 of the entire cohort (**Supplementary Figure S7**). Thus, ALLCatchR can identify a subset of
269 samples with worse performance for subtype allocation due to lower blast infiltration. Gene
270 expression profiles were also informative for patient’s sex and disease immunophenotype. To
271 enable gene expression based cross-validation of these important clinical baseline
272 characteristics, we implemented sub-classifiers to the samples immunophenotype (pro-B vs.
273 common-/pre-B ALL; accuracy of 0.871 in the validation data) and patient’s sex (accuracy:

274 0.991 in validation data set, **Figure 3C**). ALLCatchR thus provides a cross-validation of clinical
275 baseline variables and allows imputation of missing values.

276 **Shared gene expression patterns suggest distinct cells of origin for BCP-ALL subtypes**

277 The cell of origin for BCP-ALL cases remains to be defined, with immunophenotyping according
278 to EGIL criteria²⁰ representing a framework for orientation. An improved understanding of
279 underlying lymphopoiesis trajectories is especially warranted regarding current
280 immunotherapies which rely on differentiation-stage- and lineage-specific markers as
281 therapeutic targets. To map BCP-ALL subtypes to underlying B lymphopoiesis trajectories, we
282 established a reference of normal human B lymphopoiesis for 7 differentiation stages from
283 hematopoietic stem cells to mature bone marrow B cell subsets (**Figure 4A**), based on
284 established definitions²¹. Expression profiles were obtained from ultra-low input RNA-Seq of
285 FACS sorted bone marrow samples of healthy adult donors (n=4). Unsupervised analysis of
286 variable expressed genes grouped samples according to the developmental course (**Figure**
287 **4B**). Stage specific gene sets were obtained by multi-comparison ANOVA on normalized counts
288 (vst), yielding well discriminative definitions (**Figure 4C; Supplementary Table S6**). Analysis of
289 immunoglobulin rearrangements using droplet PCR indicated initiation of D_H-J_H
290 rearrangements in sorted pro-B cells while $V_H-(D)J_H$ rearrangements were first observed in pre-
291 B II Large cells and class switch recombination occurred exclusively in the most mature B cells,
292 providing an immunogenomic differentiation trajectory²² which independently confirms our
293 sorting strategy (**Supplementary Figure S8**). We implemented this newly established model
294 of human B lymphopoiesis in ALLCatchR using ssGSEA to define the proximity of each BCP-ALL
295 sample to all 7 lymphopoiesis stages (**Figure 4D; Supplementary Figure S9**). Medians of these
296 enrichment scores across samples revealed distinct patterns of enrichments suggesting
297 shared stages of origin for BCP-ALL subtypes (pro-B / pre-B I / pre-B I to pre-B II Large transition
298 / pre-B II Large; **Supplementary Figure S9**) with similar patterns in pediatric and adult data
299 sets (**Supplementary Figure S10**). Most BCP-ALL subtypes and the majority of all cases showed
300 highest similarity to the pre-B I stage (Figure 4D). *KMT2A*-rearranged and *PAX5* P80R ALL
301 showed a clearly distinct enrichment pattern favoring an earlier pro-B differentiation stage of
302 origin (**Figure 4E**). In contrast, *CEBP*, *HLF*, *IKZFN1 N159Y*, *MEF2D*, *NUTM1* and *TCF3::PBX1* were
303 grouped in a cluster with highest enrichment in transition of pre-B-I to pre-B-II large stage and
304 *BCL2/MYC* showed the highest degree of similarity exclusively to pre-B II Large differentiation
305 stage (**Figure 4D**). These observations confirm expectations for the extremes of this trajectory

306 (*KMT2A* and *BCL2/MYC*).^{23,24} A recently reported mouse model of *PAX5* P80R ALL²⁵ established
307 a pro-B differentiation arrest as initial event in *PAX5* P80R homozygous models, supporting a
308 pro-B origin of this leukemia subtype or at least an altered *PAX5* function inducing a pro-B like
309 phenotype in P80R mutated cases (**Figure 4E**). Thus, specific enrichment patterns of normal
310 lymphopoiesis are shared between molecular subtypes, suggesting distinct stages of
311 transition from normal to leukemic lymphopoiesis. We have included this model in ALLCatchR.
312 Comparison of EGIL immunophenotypes to gene-expression-defined stages of origin indicated
313 expected enrichments (pro-B stage in pro-B immunophenotype / pre-B II Large in pre-B
314 immunophenotypes; **Figure 4F**) but nearly all gene-expression-based differentiation stages
315 were represented in each immunophenotype. BCP-ALL subtypes were more closely related to
316 gene-expression-based differentiation stages as to EGIL immunophenotypes, suggesting that
317 ALLCatchR identifies developmental underpinnings of BCP-ALL drivers at higher resolution.

318 **BCP-ALL subtype-defining gene sets indicate shared signaling trajectories**

319 Definitions of BCP-ALL subtype specific gene expression signatures depend on the size and
320 composition of the remaining cohort used as comparator. We made use of the aggregated
321 transcriptome profiles of 21 BCP-ALL subtypes to define subtype specific gene expression
322 profiles based on the largest data set (n=3,532) available till date, representing different age
323 groups, cohorts, and sequencing methods. UMAP clustering of all samples according to LASSO
324 selected subtype specific gene sets indicated a clear separation of molecular subtypes
325 independently of the contributing cohorts (**Figure 5A**). To characterize subtype specific gene
326 expression profiles beyond top discriminative features, we performed differential gene
327 expression analysis for each subtype compared to the remaining cohort. A median of 673
328 differentially expressed genes per subtype were identified (range: 144– 1465; fold change:
329 <1.5-log₂-fold, FDR: <0.001; **Figure 5B**). Overlap between these gene sets was very low
330 (**Supplementary Figure S11**) indicating that subtype-specific differences are represented in
331 broad gene regulatory programs. Subtype specific gene expression profiles were provided as
332 a resource in **Supplementary Tables S7-28**. To explore the potential of this dataset to reveal
333 underlying biological functions, we performed ssGSEA for canonical signaling pathways
334 (MSigDB Hallmark / KEGG gene sets). Analysis of pathways top differentially enriched in BCP-
335 ALL subtypes (one-way ANOVA) indicated previously unrecognized clusters of subtypes with
336 enrichment in cytokine receptor / JAK-STAT signaling (Ph-pos, Ph-like, *ZNF384*, Hyperdiploid,
337 *iAMP21*) or WNT-/beta catenin/ hedgehog signaling (*ETV6::RUNX1* and -like, *CDX2/UBTF*),

338 which together represented the majority of subtypes with a putative pre-B-I cell of origin
339 (**Figure 5C**). For the remaining subtypes an enrichment in MYC-/MTOR signaling was observed
340 in subtypes of both, a more and less mature differentiation stage of origin (pro-B: *KMT2A*, *PAX*
341 *P80R* / pre-B I to pre-B II large: *BCL2/MYC*, *IKZF1 N159Y*, *MEF2D*; **Figure 5C**). Thus, enrichment
342 analysis for canonical signaling pathways independently grouped together BCP-ALL subtypes
343 form similar underlying B lymphopoiesis differentiation stages. ALLCatchR not only provides a
344 systematic gene expression analysis for accurate identification of molecular BCP-ALL subtypes
345 but also enables insights into underlying disease biology which is closely interconnected with
346 subtype nosology.

347 **Discussion**

348 Risk stratification based on molecular disease subtypes has contributed to the remarkable
349 improvement in outcomes of patients with BCP-ALL in the last decades and has provided
350 guidance for target specific treatments. Current nosology of BCP-ALL includes up to 26 specific
351 subtypes (WHO-HAEM5/ICC)^{1,2}, exceeding the capability of cytogenetic and molecular genetic
352 techniques which have so far been combined for molecular subtype allocation. Transcriptome
353 sequencing provides informative gene expression profiles and allows identification of
354 underlying driver gene fusions and more recently also driver single nucleotide variants and
355 karyotypes. Analysis of gene expression profiles for molecular subtype allocation is still not
356 standardized, despite its potential for validating genomic driver calls and for subtype
357 allocation of samples with missed genomic drivers.⁴

358 We have developed ALLCatchR, a pre-trained machine learning classifier which allows
359 molecular subtype allocation in independent hold-out data with >95% accuracy. ALLCatchR is
360 the only tool which systematically provides allocation to all gene expression defined subtypes
361 of the ICC classification, including novel *CDX2/UBTF* ALL^{4,26–28} and *CEBP/ZEB2*^{29–31}. Comparable
362 published approaches (ALLSorts, ALLspice) also achieved accurate predictions. However,
363 ALLCatchR achieved superior performance through enabling more correct subtype allocations
364 especially in a real-world adult BCP-ALL data set from a diagnostic laboratory (MLL)⁸, probably
365 due to incorporation of similar data from an independent adult cohort in the training set
366 (GMALL)^{3,4}. Immunophenotyping is a routine diagnostic in BCP-ALL and provides putative
367 differentiation stages of origin with ‘pro-B’ immunophenotype used as high-risk marker in
368 some treatment stratification systems. EGIL definitions²⁰ were derived from murine B

369 lymphopoiesis. Projecting BCP-ALL samples to our newly established reference of normal
370 lymphopoiesis yielded novel insights into differentiation stages of origin shared between BCP-
371 ALL subtypes. Interestingly, *KMT2A* and *PAX5* P80R ALL, showed a strong proximity to normal
372 pro-B cells, the most immature B lymphoid stage analyzed. These observations are in line with
373 very recent single cell analyses suggesting a pro-B or even pre-pro-B origin of *KMT2A* ALL^{24,32}
374 and murine models of *PAX5* P80R ALL showing that homozygous *PAX5* P80R induces a pro-B
375 differentiation arrest in lymphopoiesis before full transformation through acquisition of
376 additional driver events.²⁵ Here, ALLCatchR analysis based on our large aggregated reference
377 cohort confirmed these observations of smaller cohorts^{24,32}, preclinical models²⁵ and previous
378 assumptions on red-directed *PAX5* functionality in *PAX5* P80R ALL^{3,5}. Gene-expression-based
379 definitions of developmental stages in BCP-ALL were more closely related to BCP-ALL subtypes
380 than immunophenotypes, suggesting that selection for leukemogenic drivers occurs in a
381 differentiation-stage specific manner.

382 ALLCatchR is based on the largest cohort of BCP-ALL gene expression profiles across age
383 groups and molecular subtypes available till date. We make use of this aggregated data to
384 provide subtype defining gene sets for normal and leukemic B lymphopoiesis as an
385 independent research resource. Although only a small minority of samples remain
386 ‘unassigned’, novel subtype candidates are being discussed (e.g.; *IDH1/2* mutated ALL, Low
387 hyperdiploid ALL)^{5,26}. ALLCatchR is a freely available open-source tool providing a conceptual
388 and technical framework which can easily be extended for incorporation of novel subtypes
389 and additional predictive models. When combined with already established approaches for
390 calling of genomic drivers (e.g., gene fusions), ALLCatchR will complement the essential
391 prerequisites for the transition of RNA-Seq from research to routine application in clinical
392 diagnostics.

393 **Acknowledgements:**

394 This study was in part funded by the Deutsche Forschungsgemeinschaft (DFG, German
395 Research Foundation) – project number 444949889 (KFO 5010/1 Clinical Research Unit
396 ‘CATCH ALL’ to L.B., A.H., M.P.H., M.N., M.B., and C.D.B.), and project number 413490537
397 (Clinician Scientist Program in Evolutionary Medicine to B.T.H.) and Deutsche Jose Carreras

398 Leukämie Stiftung (DJCLS 01R/2016 to L.B. and C.D.B, DJCLS R 15/11 and DJCLS 06R/2019 to
399 M.Br.) and the Czech Health Research Council (NU20-07-00322 to M.Z. and J.T.)

400

401 We gratefully appreciate critical contributions from Saskia Kohlscheen and Matthias Ritgen
402 for the development of the healthy donor FACS sort panel and Monika Szczepanowski for
403 contributing to sample collection critical discussion of the manuscript. We are indebted to
404 Christian Peters and Esther Schiminsky for performing the FACS sorts.

405 **Author contributions:**

406 T.B., M.Br., C.D.B. and L.Bas. designed the study; T.B. and L.Bas. established models for
407 molecular subtype allocation and B cell developmental stages and developed the classifier;
408 B.T.H., L.Bas. and C.D.B. conceived the clinical trial to obtain healthy bone marrow samples;
409 B.T.H., E.A. and L.Bas. established the normal donor FACS panel, B.T.H. and M.Bu. performed
410 FACS sorting; T.B., B.T.H, A.M.H., N.K., L.Bar., S.B., J.K., M.B. established bioinformatic
411 workflows and performed analyses of BCP-ALL and healthy donor gene expression profiles;
412 J.Z. and Chr. K. developed and tested the CRAN package for ALLCatchR distribution, W.W.,
413 M.Z., Z.A., P.C., G.C., M.S., M.N., N.G., A.K.B., J.T., C.H. contributed BCP-ALL sequencing data
414 and validated ground truth and/or contributed to the classifier concept; L.Bas. and C.D.B.
415 supervised the project; T.B., C.D.B. and L.Bas. drafted the first version of the manuscript; all
416 authors revised and approved the final version of the manuscript.

417 **Competing Interests:**

418 The authors have no competing interests to declare.

419 **Data Availability Statement:**

420 ALLCatchR is freely available as an R-package through
421 <https://github.com/ThomasBeder/ALLCatchR>. Transcriptome sequencing data of bone
422 marrow samples from healthy donors were deposited at the European Genome Phenome
423 archive. The accession number will be provided after acceptance of the manuscript. BCP-ALL

424 transcriptome profiles haven been deposited in open or controlled access archives by the
425 authors of the original publications.

426 References:

- 427 1 Alaggio R, Amador C, Anagnostopoulos I, Attygalle AD, Araujo IB de O, Berti E *et al*. The 5th edition of the
428 World Health Organization Classification of Haematolymphoid Tumours: Lymphoid Neoplasms. *Leukemia*
429 2022; **36**: 1720–1748.
- 430 2 Arber DA, Orazi A, Hasserjian RP, Borowitz MJ, Calvo KR, Kvasnicka HM *et al*. International Consensus
431 Classification of Myeloid Neoplasms and Acute Leukemia: Integrating Morphological, Clinical, and
432 Genomic Data. *Blood* 2022; : blood.2022015850.
- 433 3 Bastian L, Schroeder MP, Eckert C, Schlee C, Sanchez JO, Kämpf S *et al*. PAX5 biallelic genomic alterations
434 define a novel subgroup of B-cell precursor acute lymphoblastic leukemia. *Leukemia* 2019; **33**: 1895–1909.
- 435 4 Bastian L, Hartmann AM, Beder T, Hänzelmann S, Kässens J, Bultmann M *et al*. UBTF::ATXN7L3 gene fusion
436 defines novel B cell precursor ALL subtype with CDX2 expression and need for intensified treatment.
437 *Leukemia* 2022. doi:10.1038/s41375-022-01557-6.
- 438 5 Gu Z, Churchman ML, Roberts KG, Moore I, Zhou X, Nakitandwe J *et al*. PAX5-driven subtypes of B-
439 progenitor acute lymphoblastic leukemia. *Nat Genet* 2019; **51**: 296–307.
- 440 6 Zaliova M, Stuchly J, Winkowska L, Musilova A, Fiser K, Slamova M *et al*. Genomic landscape of pediatric B-
441 other acute lymphoblastic leukemia in a consecutive European cohort. *Haematologica* 2019; **104**: 1396–
442 1406.
- 443 7 Chouvarine P, Antić Ž, Lentjes J, Schröder C, Alten J, Brüggemann M *et al*. Transcriptional and Mutational
444 Profiling of B-Other Acute Lymphoblastic Leukemia for Improved Diagnostics. *Cancers* 2021; **13**: 5653.
- 445 8 Walter W, Shahswar R, Stengel A, Meggendorfer M, Kern W, Haferlach T *et al*. Clinical application of
446 whole transcriptome sequencing for the classification of patients with acute lymphoblastic leukemia. *BMC*
447 *Cancer* 2021; **21**: 886.
- 448 9 D. Nicorici, M. Satalan, H. Edgren, S. Kangaspeska, A. Murumagi, O. Kallioniemi, S. Virtanen, O. Kilku.,
449 FusionCatcher – a tool for finding somatic fusion genes in paired-end RNA-sequencing data, bioRxiv, Nov.
450 2014, DOI:10.1101/011650. .
- 451 10 Rehn J, Mayoh C, Heatley SL, McClure BJ, Eadie LN, Schutz C *et al*. RaScALL: Rapid (Ra) screening (Sc) of
452 RNA-seq data for prognostically significant genomic alterations in acute lymphoblastic leukaemia (ALL).
453 *PLoS Genet* 2022; **18**: e1010300.
- 454 11 Bařinka J, Hu Z, Wang L, Wheeler DA, Rahbarinia D, McLeod C *et al*. RNAseqCNV: analysis of large-scale
455 copy number variations from RNA-seq data. *Leukemia* 2022; **36**: 1492–1498.
- 456 12 Boer JM, Marchante JRM, Evans WE, Horstmann MA, Escherich G, Pieters R *et al*. BCR-ABL1-like cases in
457 pediatric acute lymphoblastic leukemia: a comparison between DCOG/Erasmus MC and COG/St. Jude
458 signatures. *Haematologica* 2015; **100**: e354–e357.
- 459 13 Schmidt B, Brown LM, Ryland GL, Lonsdale A, Kosasih HJ, Ludlow LE *et al*. ALLSorts: an RNA-Seq subtype
460 classifier for B-cell acute lymphoblastic leukemia. *Blood Adv* 2022; **6**: 4093–4097.
- 461 14 Mäkinen V-P, Rehn J, Breen J, Yeung D, White DL. Multi-Cohort Transcriptomic Subtyping of B-Cell Acute
462 Lymphoblastic Leukemia. *Int J Mol Sci* 2022; **23**: 4574.
- 463 15 Tibshirani R. THE LASSO METHOD FOR VARIABLE SELECTION IN THE COX MODEL. *Stat Med* 1997; **16**: 385–
464 395.
- 465 16 Friedman J, Hastie T, Tibshirani R. Regularization Paths for Generalized Linear Models via Coordinate
466 Descent. *J Stat Softw* 2010; **33**: 1–22.
- 467 17 Kursu MB, Rudnicki WR. Feature Selection with the **Boruta** Package. *J Stat Softw* 2010; **36**.
468 doi:10.18637/jss.v036.i11.

- 469 18 Kuhn M. Building Predictive Models in R Using the **caret** Package. *J Stat Softw* 2008; **28**.
470 doi:10.18637/jss.v028.i05.
- 471 19 Foroutan M, Bhuvu DD, Lyu R, Horan K, Cursons J, Davis MJ. Single sample scoring of molecular
472 phenotypes. *BMC Bioinformatics* 2018; **19**: 404.
- 473 20 Bene MC, Castoldi G, Knapp W, Ludwig WD, Matutes E, Orfao A *et al*. Proposals for the immunological
474 classification of acute leukemias. European Group for the Immunological Characterization of Leukemias
475 (EGIL). *Leukemia* 1995; **9**: 1783–1786.
- 476 21 van Zelm MC, van der Burg M, de Ridder D, Barendregt BH, de Haas EFE, Reinders MJT *et al*. Ig gene
477 rearrangement steps are initiated in early human precursor B cell subsets and correlate with specific
478 transcription factor expression. *J Immunol Baltim Md 1950* 2005; **175**: 5912–5922.
- 479 22 Zoutman WH, Nell RJ, Versluis M, Pico I, Khanh Vu TH, Verdijk RM *et al*. A novel digital PCR-based method
480 to quantify (switched) B cells reveals the extent of allelic involvement in different recombination
481 processes in the IGH locus. *Mol Immunol* 2022; **145**: 109–123.
- 482 23 Wagener R, López C, Kleinheinz K, Bausinger J, Aukema SM, Nagel I *et al*. IG-MYC+ neoplasms with
483 precursor B-cell phenotype are molecularly distinct from Burkitt lymphomas. *Blood* 2018; **132**: 2280–
484 2285.
- 485 24 Chen C, Yu W, Alikarami F, Qiu Q, Chen C-H, Flournoy J *et al*. Single-cell multiomics reveals increased
486 plasticity, resistant populations, and stem-cell-like blasts in KMT2A-rearranged leukemia. *Blood* 2022; **139**:
487 2198–2211.
- 488 25 Jia Z, Hu Z, Damirchi B, Han TT, Gu Z. Characterization of PAX5 Mutations in B Progenitor Acute
489 Lymphoblastic Leukemia. *Blood* 2022; **140**: 1001–1002.
- 490 26 Yasuda T, Sanada M, Kawazu M, Kojima S, Tszuzuki S, Ueno H *et al*. Two novel high-risk adult B-cell acute
491 lymphoblastic leukemia subtypes with high expression of *CDX2* and *IDH1/2* mutations. *Blood* 2022; **139**:
492 1850–1862.
- 493 27 Passet M, Kim R, Gachet S, Sigaux F, Chaumeil J, Galland A *et al*. Concurrent *CDX2* cis-deregulation and
494 *UBTF-ATXN7L3* fusion define a novel high-risk subtype of B-cell ALL. *Blood* 2022; : blood.2021014723.
- 495 28 Kimura S, Montefiori L, Iacobucci I, Zhao Y, Gao Q, Paietta EM *et al*. Enhancer retargeting of *CDX2* and
496 *UBTF::ATXN7L3* define a subtype of high-risk B-progenitor acute lymphoblastic leukemia. *Blood* 2022; :
497 blood.2022015444.
- 498 29 Li J-F, Dai Y-T, Lilljebjörn H, Shen S-H, Cui B-W, Bai L *et al*. Transcriptional landscape of B cell precursor
499 acute lymphoblastic leukemia based on an international study of 1,223 cases. *Proc Natl Acad Sci U S A*
500 2018; **115**: E11711–E11720.
- 501 30 Akasaka T, Balasas T, Russell LJ, Sugimoto K, Majid A, Walewska R *et al*. Five members of the CEBP
502 transcription factor family are targeted by recurrent IGH translocations in B-cell precursor acute
503 lymphoblastic leukemia (BCP-ALL). *Blood* 2007; **109**: 3451–3461.
- 504 31 Zaliova M, Potuckova E, Lukes J, Winkowska L, Starkova J, Janotova I *et al*. Frequency and prognostic
505 impact of ZEB2 H1038 and 1072 mutations in childhood B-other acute lymphoblastic leukemia.
506 *Haematologica* 2020. doi:10.3324/haematol.2020.249094.
- 507 32 Khabirova E, Jardine L, Coorens THH, Webb S, Treger TD, Engelbert J *et al*. Single-cell transcriptomics
508 reveals a distinct developmental state of KMT2A-rearranged infant B-cell acute lymphoblastic leukemia.
509 *Nat Med* 2022; **28**: 743–751.
- 510 33 Leek JT, Johnson WE, Parker HS, Jaffe AE, Storey JD. The sva package for removing batch effects and other
511 unwanted variation in high-throughput experiments. *Bioinformatics* 2012; **28**: 882–883.
- 512 34 Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with
513 DESeq2. *Genome Biol* 2014; **15**: 550.

514 Figure Legends

515 **Figure 1. ALLCatchR predicts molecular BCP-ALL subtypes based on gene expression count**
516 **data with high accuracy. (A)** Heatmap showing the prediction scores for 21 gene expression
517 defined BCP-molecular subtypes (WHO-HAEM5 / ICC) in n=3,308 samples of the entire BCP-
518 ALL cohort (after removal of duplicate samples and samples with two primary subtype
519 allocations; n=217) samples. Molecular subtypes had been defined in the six original studies
520 (GMALL, St Jude, CLIP, MLL, MHH and RCH/PM) based on genomic driver aberrations and
521 corresponding gene expression signatures in n=2,887 cases (ground truth). Remaining cases
522 were deemed ‘unassigned’ or ‘B-other’. ALLCatchR scores are shown for the combined data
523 set of training and hold-out cohorts. **(B)** Cutoffs were defined for each BCP-ALL subtype based
524 on distribution of all ALLCatchR scores in every subtype. Cutoffs for ‘high confidence
525 predictions’ were defined to include >90% of all samples allocated to these subtypes in the
526 original data set, resulting in 0.989 accuracy of these predictions in independent hold-out data
527 set. Cutoffs for ‘Candidate predictions’ were defined to reliably exclude samples from other
528 subtypes, providing a reliable orientation for further validation of subtype assignment based
529 on genomic drivers (accuracy: 0.851). ‘Low-confidence’ predictions indicate samples from
530 different subtypes or samples where subtype allocation cannot be performed. These were
531 considered ‘unclassified’ for further analysis. **(C)** The proportions of confidence categories for
532 true and false predictions in the training and hold-out data sets are shown. A prediction was
533 considered ‘true’ if the sample received the same subtype allocation as in the original study.
534 ‘False’ predictions represent allocations to other subtypes than the subtype assigned in the
535 original study. For comparison, ‘unassigned’ / ‘B-other’ samples from the holdout data sets
536 are shown. **(D)** Confusion matrices relate ALLCatchR predictions to the ground truth in training
537 samples (left) and holdout cohorts (right). By design, the training cohort did not contain
538 ‘unassigned’ / ‘B-other’ samples. In the hold-out data, n=111 samples had been defined as
539 ‘unassigned’ / ‘B-other’ and predictions for these are also shown. Supplementary Figure S5A
540 and Supplementary Table S5 indicate how ALLCatchR predictions in ‘unassigned / B-other’
541 samples are supported by corresponding genomic drivers in 72.1% of ‘high confidence’ and
542 27.1% of ‘candidate’ predictions.

543

544 **Figure 2. ALLCatchR accuracy for subtype allocation is consistently high across cohorts and**
545 **BCP-ALL subtypes. (A)** Sankey diagrams indicate ALLCatchR subtype allocations and
546 corresponding subtype validated ground truth in the training cohort and the individual

547 holdout data sets. ‘Acc.’ Indicated accuracy in the corresponding data set. **(B)** Bar charts
548 indicated sensitivity and specificity for the individual subtypes in the training and hold-out
549 data. Validated ground truth was used to define true positive cases, i.e. belonging to this
550 subtype and true negative cases, i.e. not belonging to this subtype. Values were obtained as
551 fraction of true positive cases from all cases defined by ALLCatchR as belonging to this subtype
552 (sensitivity) and as fraction of true negative cases from all cases defined by ALLCatchR as not
553 belonging to this subtype (specificity).

554

555 **Figure 3. ALLCatchR predicts sample blast counts, patient’s sex and immunophenotype**
556 **based on gene expression data. (A)** For GMALL (n=302), MLL (n=282) and RCH/PM (n=77)
557 sample blast counts obtained by cytology or flow cytometry were available. GMALL and MLL
558 cohorts were separately used for training two classifiers in a 10-fold cross-validation scheme
559 with the same machine learning algorithms used for subtype prediction. GMALL and MLL
560 classifiers were validated on each other, and both were validated on the RCH/PM data. Best
561 performing methods in terms of the Root Mean Squared Error (RSME) on the training data are
562 shown. Training two classifiers on independent data sets allowed for the validation on each
563 other and both were combined for final predictions. Blast count predictions had a good
564 correlation to measured counts i.e., $\rho=0.590$ in GMALL and $\rho=0.771$ in MLL. Moreover,
565 predicting MLL samples with the classifier trained on GMALL achieved a similar performance
566 as the classifier trained on MLL samples and *vice versa*. **(B)** Since both, GMALL and MLL
567 classifiers had a good performance and were generalizable, predictions from both are
568 combined in ALLCatchR. **(C)** Sub-classifiers for immunophenotype and patient’s sex were
569 developed using SVMlinear and ranger machine learning models respectively. An
570 immunophenotype classifier was trained on GMALL samples (n=413 common-B / pre-B and
571 n=66 pro-B) and validated on MLL data (n=168 common-B / pre-B and n=64 pro-B) with
572 available EGIL immunophenotypes. A patient sex classifier was trained on n=357 GMALL
573 samples (female=165, male=192) analogous to the subtype classifier. For validation n=1892
574 St Jude samples with known sex (female=850, male=1042) were used. Corresponding
575 accuracies, sensitivities and specificities are shown for these sub-classifiers.

576

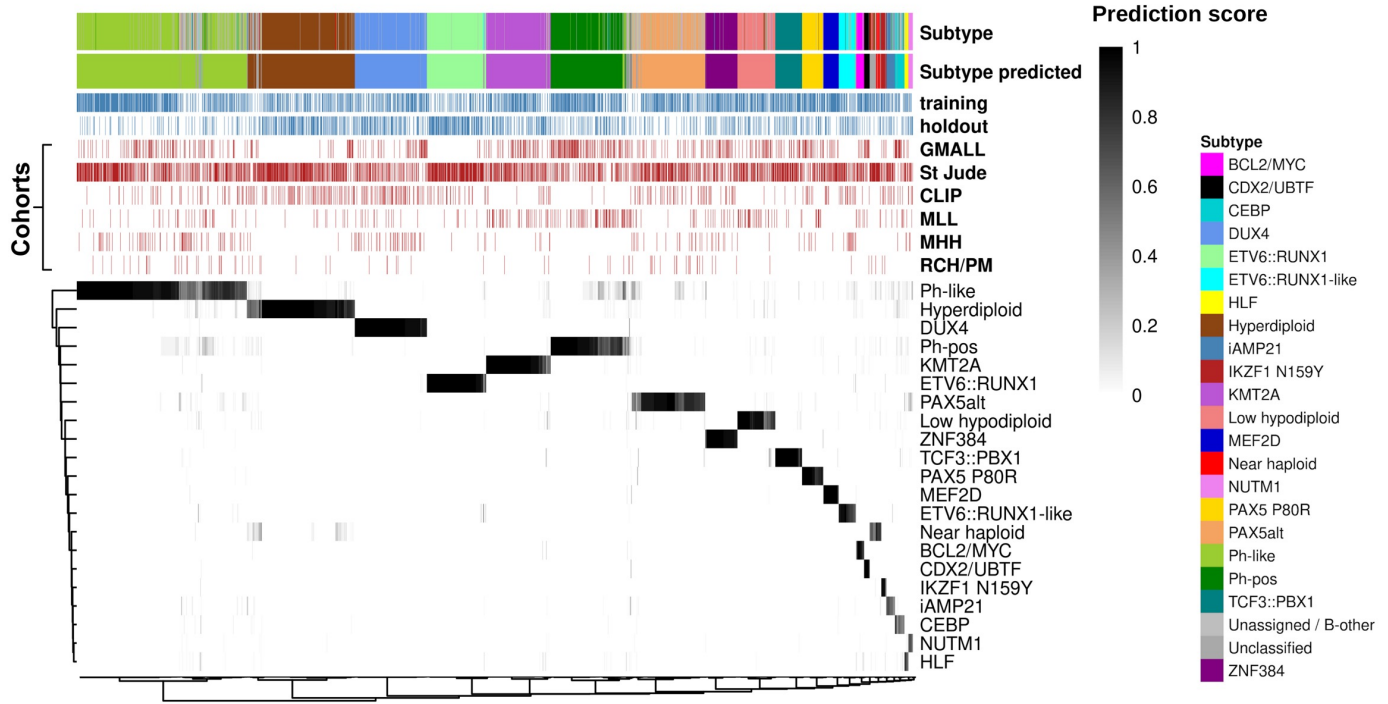
577 **Figure 4. ALLCatchR identifies B cell developmental trajectories underlying BCP-ALL**
578 **subtypes. (A)** To establish a reference map of human B lymphopoiesis, we obtained bone
579 marrow samples from healthy adult donors (n=4) and used a 9-color antibody panel for FACS

580 sorting of 7 B lymphopoiesis stages following described definitions²¹ after pre-enrichment of
581 wanted populations. Lin⁻ selection included CD3, CD33, CD56, CD14, CD66c, CD138.
582 Antibodies used are shown in Supplementary Table S4. Supplementary Figure S8 shows
583 immunogenomic profiling of immune gene rearrangements in support of the applied sorting
584 strategy. **(B)** Ultra-low input RNA-Seq was performed for total RNA to obtain stage-specific
585 gene expression. Uniform manifold approximation plot (UMAP) shows clustering of human B
586 lymphopoiesis stages based on 400 most variable expressed genes. **(C)** Multi comparison
587 ANOVA on normalized (vst) count data was performed to obtain differentiation-stage specific
588 gene sets. Heatmap depicts single sample gene set enrichment analyses (singscore)¹⁹ of B
589 lymphopoiesis subsets (columns) to stage defining gene sets (rows). **(D)** BCP-ALL samples with
590 known subtype allocation (n=2,887) were used for single sample gene set enrichment analysis
591 with B lymphopoiesis-specific gene sets obtained from (C). Supplementary Figure S9 shows
592 enrichment patterns of individual samples from all BCP-ALL subtypes for all differentiation
593 stages. Heatmap depicts averaged enrichment scores for all BCP-ALL subtypes and all B
594 lymphopoiesis stages grouped by unsupervised clustering. Normal progenitors with closest
595 proximity to BCP-ALL subtypes representing putative cells-of-origin are annotated on top.
596 Supplementary Figure S10 provides separate analyses for pediatric and adult patients
597 indicating a high degree of similarity. **(E)** *KMT2A* rearranged and PAX5 P80R ALL had both the
598 highest enrichment towards pro-B supporting a shared developmental origin (also depicted in
599 Supplementary Figure S9). **(F)** Comparison of gene expression defined differentiation stages
600 and EGIL immunophenotypes are shown for n=711 samples with available gene expression
601 data.

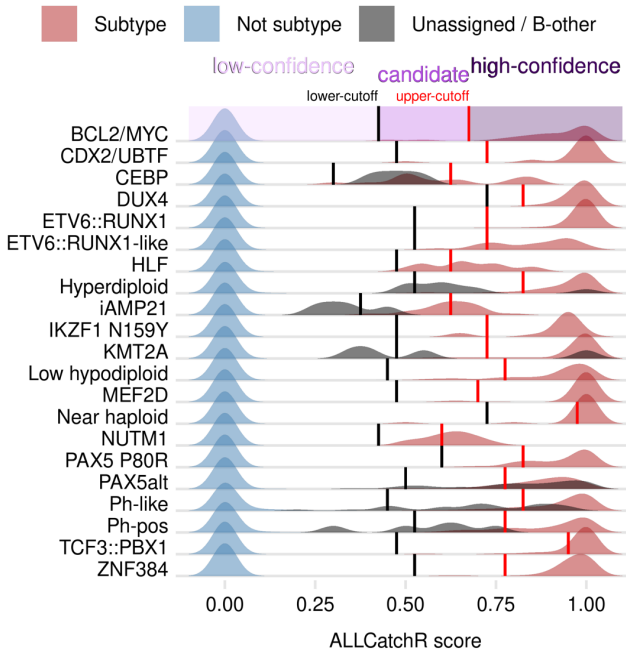
602
603 **Figure 5. The gene expression landscape in BCP-ALL. (A)** UMAP plot showing all n=3,308
604 samples used in this study. Count data from the six data sets was batch corrected using the
605 *sva* package³³ and TPM values calculated. The plot is based on 2,802 genes selected by LASSO
606 for training of ALLCatchR. Cohorts are highlighted as shape. **(B).** ALLCatchR predictions were
607 used to define samples which best represented their respective molecular subtype. A total of
608 n=20 top ranking samples per subtype (exceptions with lesser samples available: HLF n=14,
609 *CEBP* n=16, *NUTM1* n=17, *IKZF1 N159Y* n=18) were used to obtain a homogenous data set
610 representing all 21 BCP-ALL subtypes (n=405). Differential gene expression analyses for each
611 subtype versus the remaining cohort using DESeq2³⁴ revealed 5,110 differentially expressed

612 genes (cutoff: 1.5- \log_2 -fold change, FDR: 0.001) used for unsupervised clustering. Color legend
613 for BCP-ALL subtypes is the same as in (A). Supplementary Figure S11 and Supplementary
614 Tables S7-S28 provide detailed information on the derived gene sets. **(C)** Canonical signaling
615 pathways (KEGG, HALLMARK gene sets; MSigDB) were used for single sample gene set
616 enrichment analysis using the BCP-ALL subcohort from (B) for balanced representation of all
617 subtypes. Enrichment scores for top variable enriched pathways are shown.

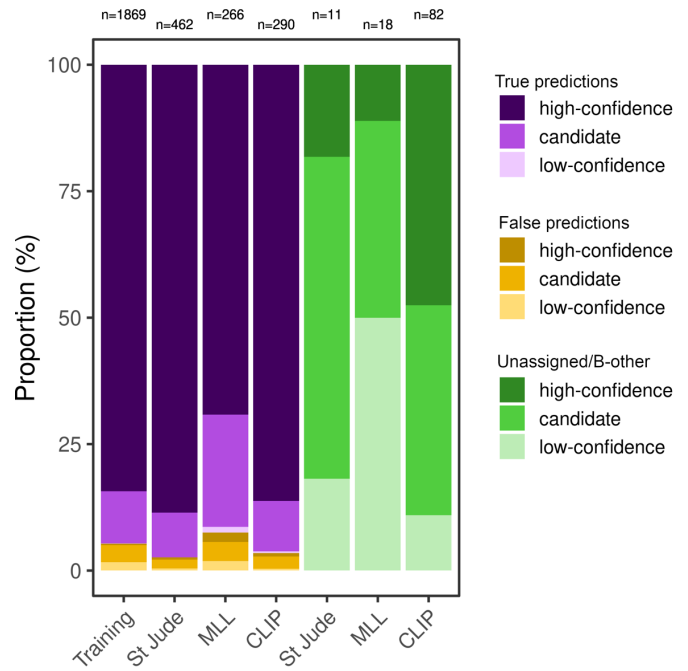
A



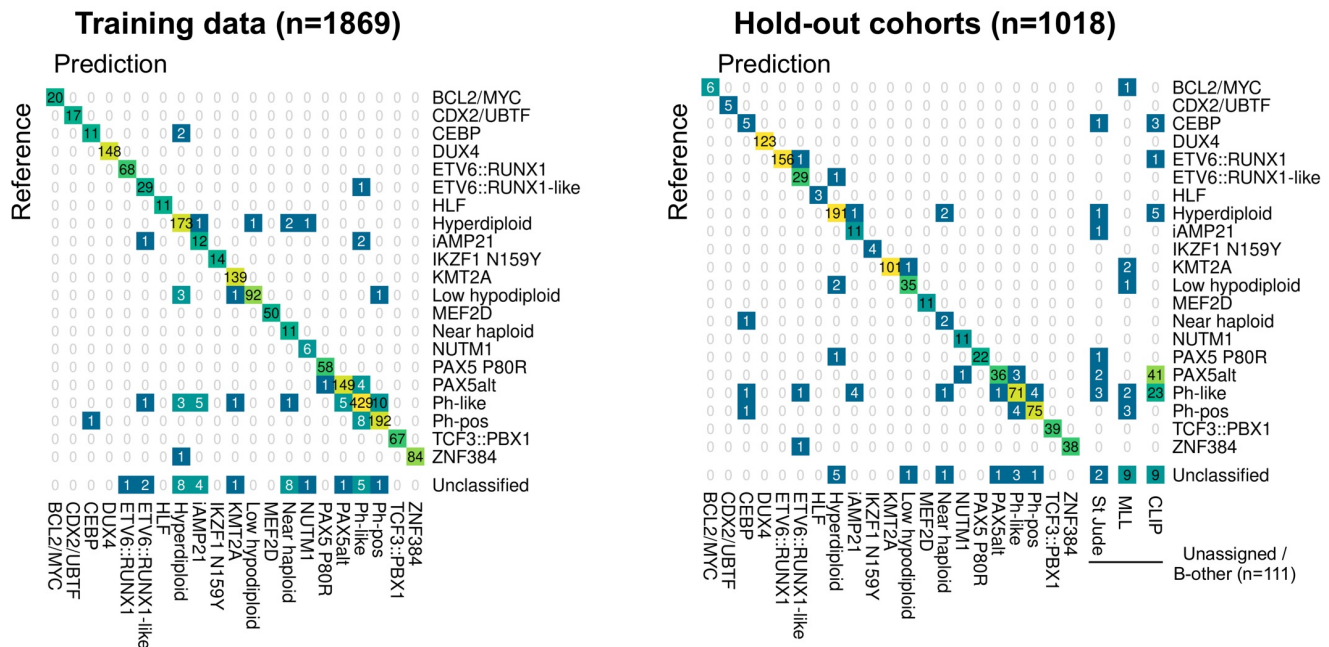
B



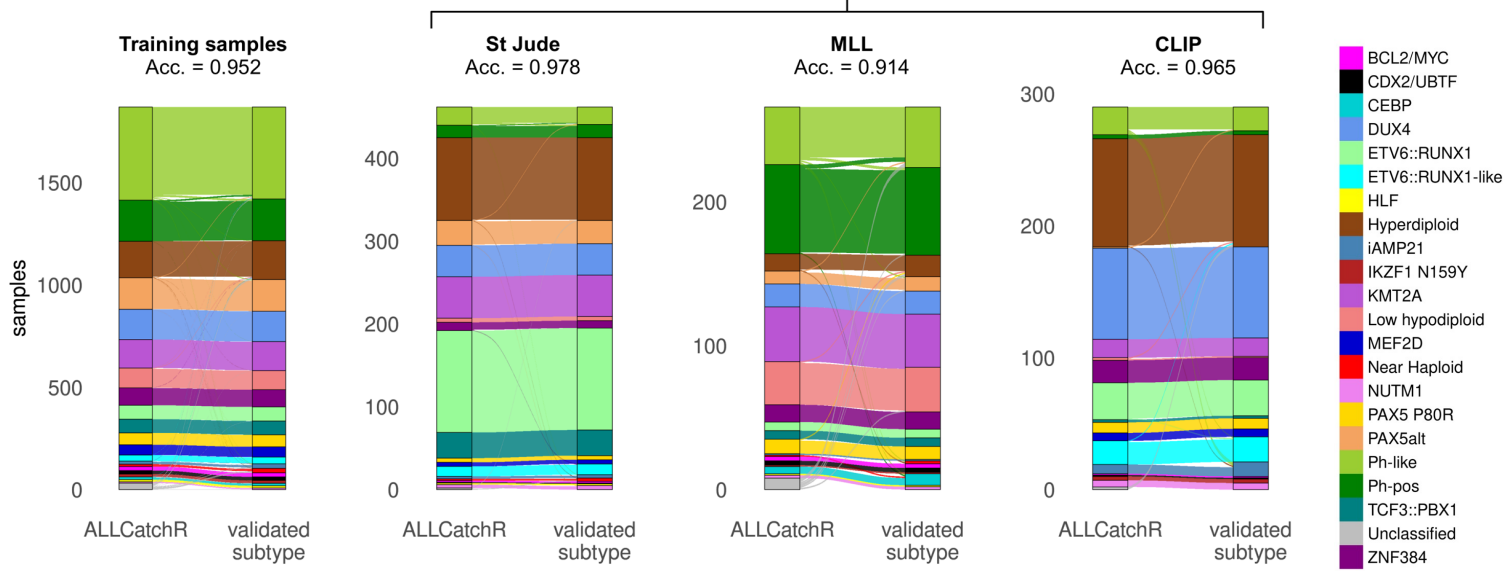
C



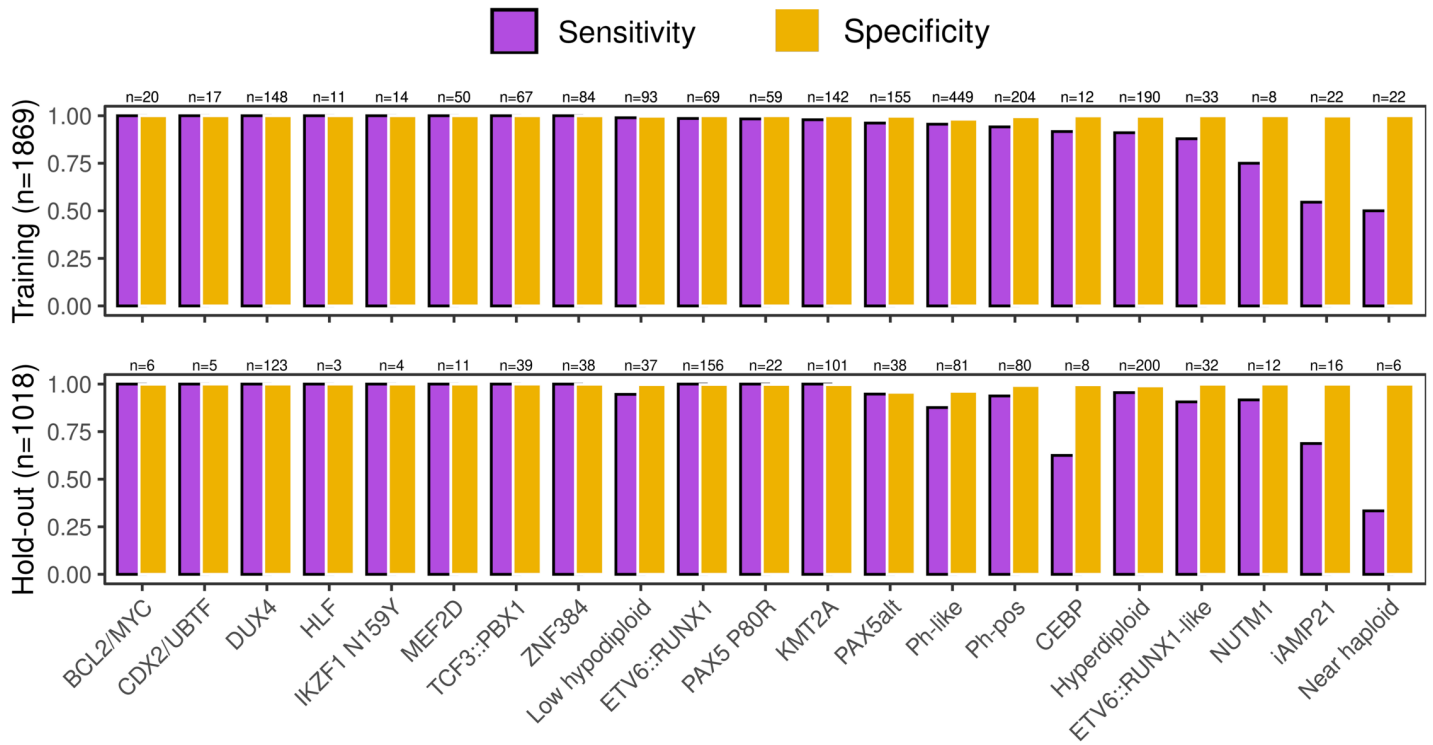
D



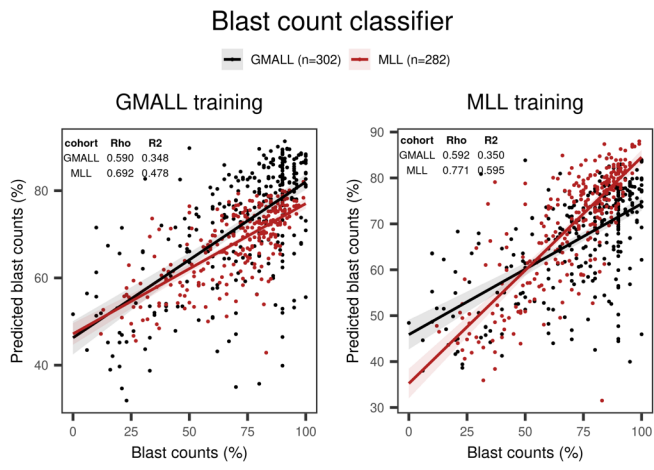
A



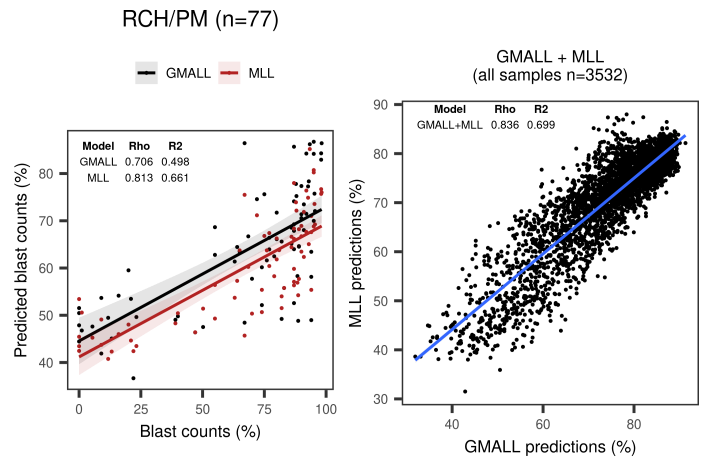
B



A



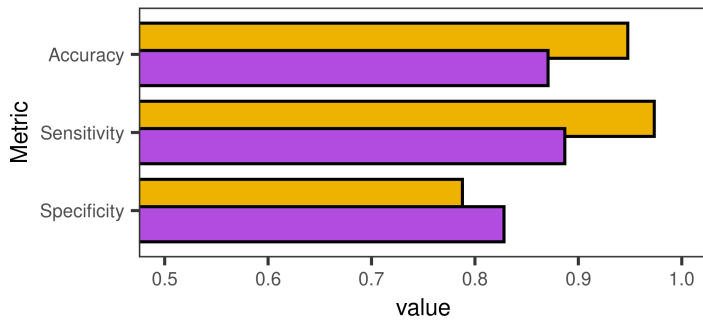
Combine GMALL and MLL blast count predictions



C

Immunophenotype classifier (common-/pre-B ALL vs. pro-B ALL)

■ Training (n=479) ■ Validation (n=232)



Patient's sex classifier

■ Training (n=357) ■ Validation (n=1892)

