

1. Supplementary Results

1.1) Simulation study to find the required sample size and the optimized the parameter α

Theoretical Simulation – We performed two simulation studies to evaluate gene-set analysis methods in the context of statistical power and type-I error, similar to a prior study [1]. The purpose of the first study was to find a minimal sample size to guarantee statistical power for gene-set signature identification, whereas the second was to decide the optimal parameter α for the FAIME algorithm and reveal the effect of gene-set size when detecting gene-set signatures. We used four different simulation scenarios for each study, considering strong and weak signal-to-noise ratios (1 and 0.5 respectively) and altered fractions of differentially expressed (DE) genes (50% and 80% respectively).

In the first simulation study, we created mimic datasets with increased sample sizes and two gene-sets (GSs) each of a fixed GS size (30 out of 5000 genes), where one GS was differentially expressed and the other was not. Larger sample size noticeably increased statistical power but not type-I error. Using 20 samples per group, we could identify GS signatures in strong signal-to-noise-ratio scenarios (**Fig. S1-C3, D3**, 60-95% statistical power), using either FAIME. α , GSEA, or GSVA. Additionally, all methods provided effective control of the type-I error rate when GS having more than 30 gene members. GSEA and GSVA slightly outperformed FAIME.1 in scenarios with weak signal-to-noise-ratio (**Fig. S1-A, B**), while all three algorithms performed equally in the scenarios with strong signal-to-noise-ratio (**Fig. S1-C, D**). We used the Bioconductor packages PGSEA and GSVA to apply the GSEA and GSVA algorithms to all gene-sets with five or more genes from the MSigDB.

In the second simulation study, we repeated all of the aforementioned steps but tested on simulated GSs of different gene sizes ($x=10, 20, 80, \text{ and } 100$). For small GSs (10-20 genes), FAIME showed statistical advantages with consistently low type-I error (**Fig. S1- a1, b1, b2**). For other GSs (≥ 30 genes), the larger a GS size the lower a statistical power derived from the same sample size in all simulation scenarios, suggests that 40-60 samples per group are required to test larger GS (>100 genes).

Alternatively, applying a larger parameter α to FAIME will benefit the analysis of small GS via weighting more sharply towards the leading expression ranks (**Equation 1**). In both simulation studies, FAIME with $\alpha =1$ by default worked more accurately than $\alpha =5$ or 10, whereas the latter effectively controlled the type-I error for small GSs.

1.2) Summary of published LSC-associated gene signatures

Table S1 summarizes nine multiple-gene studies pertaining to LSCs, three of which the original authors checked prognosis in primary AML samples. These gene signatures can be divided into the three categories below. For the comparison between AML LSC+ and normal HSC+, we excluded those “HSC” samples from **Table S1** with detectable expression of only CD34+ (neither the mature blood lineages nor

1 their committed progenitors markers) because normal CD34⁺ mononuclear cells contain hematogones (B
2 lymphocyte precursors) and CD34⁺ megakaryocytes.

3 **Stemness** - We obtained an union set of three published multi-gene signatures, which we refer to as
4 “LSCvLPC” (**Fig. S2**). Three studies respectively identified a gene-list that significantly distinguishes
5 LSC⁺ populations (CD34⁺CD38⁻) from more mature fraction of the leukemia clone, the leukemia
6 progenitor cell (LPC) populations (CD34⁺CD38⁺), purified from the same AML samples [2-4]. Ishikawa *et*
7 *al* demonstrated that such LSC⁺ exclusively recapitulates AML and retains self-renewal capacity *in vivo*
8 [2]. Gentles *et al* showed that expression of their gene-signature in bulk primary AML tumor samples was
9 associated with clinical outcomes in four independent patient cohorts (n = 1047) [4]. Unfortunately, few
10 overlaps exist among these three signatures. Only the cell surface marker CD38 was expressed lower in
11 LSC than in LPC (**Fig. S2A1**). Additionally, a potential non-stem-cell signature was involved, as both
12 studies ignored the fact that the CD34⁺CD38⁺ subpopulation may also resides LSCs [5-8]. However,
13 these three studies are more biologically sensitive by comparing pairwise cell sub-populations purified
14 from the same patients. Therefore by joining these three published gene signatures, we could generally
15 characterize leukemia stemness.

16 An improved study by Eppert *et al* verified four divergent LSC⁺ fractions using xenograft models [5].
17 They identified multiple LSC⁺ fractions in AML samples using a sensitive xenograft assay and then
18 identified a LSC-specific signature more highly expressed in LSC⁺ than in LSC⁻. They showed this
19 signature to be a significant and independent predictor of patient survival (n=445). They also identified a
20 signature specific to normal HSC⁺ fractions but not normal mature fractions and showed its prognosis in
21 primary AML [5].

22 **Malignancy** - We joined two published gene signatures which we refer to as “LSCvNHSC”. Comparing
23 refined LSC⁺ populations (Lin⁻CD34⁺CD38⁻CD90⁻) with normal hematopoietic stem cell enriched (HSC⁺)
24 populations (Lin⁻CD34⁺CD38⁻CD90⁺), Majeti *et al* identified a LSC specific signature [9]. Similarly, Saito *et*
25 *al* identified genes with significantly higher expression in AML LSC⁺ (CD34⁺CD38⁻) than in normal HSC⁺
26 (CD34⁺CD38⁻) [10].

27 **Other dataset** - De Jonge *et al* reported 50 genes that specifically high-expressed in AML CD34⁺ but not
28 AML CD34⁻ fractions when compared with normal CD34⁺ compartment [11]. Based on the summed
29 expression level of three out of the 50 genes, they suggested that a high transcript level of CD34⁺ cells
30 was associated with significant unfavorable overall survivals in two independent cohorts (n=381) of
31 normal karyotype AML. However, the statistical significance is mild as it can be achieved after
32 trichotomizing bulk samples rather than dichotomizing.

1 2. Supplementary Methods

2 2.1) Data process

3 **Gene expression data.** We downloaded the normalized gene expression values and transformed the
 4 values to a logarithmic scale (log2) when required. One RNA-seq dataset of normal HSC samples was
 5 provided by the authors [12]. The measurements were scaled in Reads Per Kilobase of exon model per
 6 Million mapped reads (RPKM) format [13], and genes were annotated by Ensemble IDs.

7 MSigDB employs only the official gene symbols and Entrez IDs, leading us to use the Bioconductor
 8 package *biomaRt* (version 2.16.0) to map all probes on a microarray to Entrez gene IDs as well as all
 9 Ensemble IDs of RNAseq to official gene symbols. To get the best coverage for custom microarray
 10 platforms, we also incorporated the corresponding custom annotation files downloaded from GEO.

11 **Functional gene-sets** - We downloaded three categories of previously defined gene-sets from Molecular
 12 Signature Database (MSigDB, version 4.0) [14]: canonical representations of biological processes
 13 compiled by domain experts (from BIOCARTA, KEGG, and REACTOME) (N=1320), gene-sets
 14 representing expression signatures of genetic and chemical perturbations (CGP, N=3402), and
 15 transcription factor or microRNA targets based on conserved cis-regulatory motifs from a comparative
 16 analysis of the human, mouse, rat, and dog genomes (N=836). The average sizes of gene-members in
 17 these three gene-set categories vary from 29 to 233. CGP gene-sets have on average 45 gene-members
 18 per set (range of 5 to 1972).

19 2.2) Hypergeometric probability analysis

20 To estimate the probability of observing n=6 overlapping genes among three instances of random
 21 sampling (**Table S2**, three DNM GSs using FAIME.5 profiles), we did the following modeling: out of a
 22 space of N=22000 genes, we randomly picked a group of A=11 genes, recorded them and put them back,
 23 then repeated for a group of B=6 genes and C=19 genes. We wished to compute the probability that out
 24 of the 36=A+B+C recorded genes, exactly n=6 genes appear twice and exactly A+B+C-2n=24 genes
 25 appear once.

26 Note that for a gene to appear exactly twice, it must appear in two of the three groups and not in the
 27 third. In particular, any number $0 \leq n-k \leq n$ of the repeated genes may have one of its groups be C, which
 28 means the remaining k genes must be in both A and B. We can arbitrarily pick A, and given those 11

29 distinct genes, there are a total of $\binom{N}{B}\binom{N}{C}$ ways to pick B and C. We can pick the k genes from A that will

30 be repeated in B in $\binom{A}{k}$ ways, and we can pick the remaining B-k genes in B in $\binom{N-A}{B-k}$ ways.

31

1 There are a total of $A+B-k$ genes in the union of A and B, $A+B-2k=17-2k$ of which only appeared
 2 once. Since we didn't observe any gene appearing in all three of the groups but we needed to have $n-k$
 3 more genes repeated for a total of n repeated genes, we chose the $n-k$ genes from the non-repeated pool

4 of $17-2k$ genes in A and B to also appear in C. This could be done in $\binom{A+B-2k}{n-k}$ ways, and we could
 5 pick the remaining $C-(n-k)=13+k$ genes in C in $\binom{N-(A+B-k)}{C-(n-k)}$ ways.

6
 7 Summing over all k's, we have our desired probability that out of the $(A+B+C)=36$ recorded genes,
 8 exactly $(A+B+C-2n)=24$ genes appear once and exactly $n=6$ genes appear twice is:

$$p(x = n | A, B, C, N) = \sum_{k=0}^n \frac{\binom{A}{k} \binom{N-A}{B-k}}{\binom{N}{B}} \frac{\binom{A+B-2k}{n-k} \binom{N-A-B+k}{C-n+k}}{\binom{N}{C}}$$

9

10 2.3) Gene Ontology semantic similarity evaluation

11 The Gene Ontology (GO) semantic similarity between pair-wise gene members of the identified 25- or 30-
 12 gene signature was estimated using Lin's method [15]. Given two genes x and y , Lin's method assigns a
 13 semantic similarity score, $Sim(x, y) = 2 \frac{sim(x, y)}{sim(x, x) + sim(y, y)}$. We employed the Bioconductor package

14 *GOSim* to run the calculation (similarity="funSimMax", similarityTerm="relevance", normalization=TRUE)
 15 [16], respecting the GO biological process and GO molecular function respectively. To estimate the
 16 empirical p-value of an observed similarity score, we ran the same calculation for 1000 pairs of randomly
 17 selected genes.

18 2.4) Enrichment analysis (EA) on the published gene lists

19 For LSC stemness (LSC+ compared to LSC-), we focused on the gene-sets that were identified by two
 20 out of the three gene lists: the joint gene-set LSCvLPC, AML stemness, and normal stemness in Eppert
 21 study (**Fig. S2A2**, yellow circles). For LSC malignancy (AML LSC+ compared to normal HSC+), we
 22 focused on the gene-sets identified by both the jointed LSCvNHSC list (**Fig. S2A2**, dash-lined pink circles)
 23 and LSC highly expressed gene-sets in De Jonge's study (**Fig. S2A2**, blue circles). To interrogate the
 24 LSC stemness, we reported the gene-sets that were significantly enriched in two out of the three gene-
 25 lists (**Fig. S2A2**, yellow circles). Specifically for enriched gene-set in Eppert study (GES30377), we
 26 merged our EA identification with three additional author-reported LSC-associated gene-sets
 27 (BENPORATH PRC2 TARGETS, PARK HSC VS MULTIPOTENT PROGENITORS UP, and IVANOVA
 28 HEMATOPOIESIS EARLY PROGENITOR). We did so because these three gene-sets are LSC-specific

1 (LSC $p < 0.05$ but HSC $p > 0.05$, see the Table S14 in the original publication) and can be mapped into the
2 newest MsigDB v4.1 version. To interrogate the malignancy, we reported the gene-sets that were
3 significantly enriched in both two gene-lists (**Fig. S2A2**, blue and pink circles). Note that there are no
4 overlaps between stemness and malignancy for the LSC+ population, suggesting their distinct properties.

5 As we observed before [17], functional-level EA exhibits remarkably greater cross-study
6 reproducibility than gene-level significance analysis. We observed 60% repetition of enriched canonical
7 pathways on average (5-100%, FET $p < 0.05$, gene count ≥ 3 , **Fig. S2B**) and 56% repetition of targets of
8 chemical and genetic perturbations (15%-73%).

9 **3. Limitation**

10 Our study has some limitations. First, FAIME is designed to identify up-regulated or down-regulated gene-
11 sets. However, genes in the same pathway are not always differentially expressed in the same direction.
12 Some disease condition associated pathways may contain both up- and down-expressed genes caused
13 by feedback loops, such as the p53 pathway [18]. Whereas GSVA and GSAA [19] can identify this type of
14 concerted gene-set expression using a non-parametric KS-test. Second, based on an arbitrary cutoff of
15 significance at the gene-set level, identifications that have borderline differential activities or modest effect
16 size may be lost when applying inter-group comparison for both FAIME and GSVA methods. On the other
17 hand, some false positive identifications met the criteria of significance, suggesting that we should apply
18 FAIME/GSEA on the data with symmetric differential expression background at gene level. Third, the
19 'inter-dataset' normalization is a straightforward use of Z-scores. This type of standardization has been
20 successfully applied to integrate gene-set scores of differentially expressed genes and of trait-associated
21 genetic markers [19]. Even so, a standard deviation parameter for the normalization of all gene-sets,
22 including over- and under-represented gene-sets, may reduce the over-representation of some gene-sets
23 while increasing the under-representation of others. In such cases, distinct standard deviation parameters
24 for the over- and under-represented gene-sets are suggested for future discussion. Finally, gene-sets
25 only reflect an approximation of biological functions or pathways and are pre-defined. Only a subset of
26 genes within a set may contribute to a gene-set expression signature. Different gene-sets may have
27 similar signatures pertaining to the same phenotype, owing to either an overlap between the gene-sets or
28 co-regulation of non-overlapping gene-sets. Grouping correlated gene-sets and extending the interaction
29 of their gene members, perhaps by modeling on additional information, is a potentially promising
30 approach to define new functional gene-sets.

31 **Reference for Supplementary Material:**

- 32 1. Hanzelmann S, Castelo R, Guinney J: **GSVA: gene set variation analysis for microarray and RNA-**
33 **seq data.** *BMC Bioinformatics* 2013, **14**:7.
- 34 2. Ishikawa F, Yoshida S, Saito Y, Hijikata A, Kitamura H, Tanaka S, Nakamura R, Tanaka T,
35 Tomiyama H, Saito N *et al*: **Chemotherapy-resistant human AML stem cells home to and**
36 **engraft within the bone-marrow endosteal region.** *Nat Biotechnol* 2007, **25**(11):1315-1321.

- 1 3. Gal H, Amariglio N, Trakhtenbrot L, Jacob-Hirsh J, Margalit O, Avigdor A, Nagler A, Tavor S, Ein-
2 Dor L, Lapidot T *et al*: **Gene expression profiles of AML derived stem cells; similarity to**
3 **hematopoietic stem cells.** *Leukemia : official journal of the Leukemia Society of America,*
4 *Leukemia Research Fund, UK* 2006, **20**(12):2147-2154.
- 5 4. Gentles AJ, Plevritis SK, Majeti R, Alizadeh AA: **Association of a leukemic stem cell gene**
6 **expression signature with clinical outcomes in acute myeloid leukemia.** *JAMA* 2010,
7 **304**(24):2706-2715.
- 8 5. Eppert K, Takenaka K, Lechman ER, Waldron L, Nilsson B, van Galen P, Metzeler KH, Poepl A,
9 Ling V, Beyene J *et al*: **Stem cell gene expression programs influence clinical outcome in human**
10 **leukemia.** *Nat Med* 2011, **17**(9):1086-1093.
- 11 6. Taussig DC, Miraki-Moud F, Anjos-Afonso F, Pearce DJ, Allen K, Ridler C, Lillington D, Oakervee H,
12 Cavenagh J, Agrawal SG *et al*: **Anti-CD38 antibody-mediated clearance of human repopulating**
13 **cells masks the heterogeneity of leukemia-initiating cells.** *Blood* 2008, **112**(3):568-575.
- 14 7. Yoshimoto G, Miyamoto T, Jabbarzadeh-Tabrizi S, Iino T, Rocnik JL, Kikushige Y, Mori Y, Shima T,
15 Iwasaki H, Takenaka K *et al*: **FLT3-ITD up-regulates MCL-1 to promote survival of stem cells in**
16 **acute myeloid leukemia via FLT3-ITD-specific STAT5 activation.** *Blood* 2009, **114**(24):5034-5043.
- 17 8. Goardon N, Marchi E, Atzberger A, Quek L, Schuh A, Soneji S, Woll P, Mead A, Alford KA, Rout R
18 *et al*: **Coexistence of LMPP-like and GMP-like leukemia stem cells in acute myeloid leukemia.**
19 *Cancer Cell* 2011, **19**(1):138-152.
- 20 9. Majeti R, Becker MW, Tian Q, Lee TL, Yan X, Liu R, Chiang JH, Hood L, Clarke MF, Weissman IL:
21 **Dysregulated gene expression networks in human acute myelogenous leukemia stem cells.**
22 *Proc Natl Acad Sci U S A* 2009, **106**(9):3396-3401.
- 23 10. Saito Y, Kitamura H, Hijikata A, Tomizawa-Murasawa M, Tanaka S, Takagi S, Uchida N, Suzuki N,
24 Sone A, Najima Y *et al*: **Identification of therapeutic targets for quiescent, chemotherapy-**
25 **resistant human leukemia stem cells.** *Science translational medicine* 2010, **2**(17):17ra19.
- 26 11. de Jonge HJ, Woolthuis CM, Vos AZ, Mulder A, van den Berg E, Kluin PM, van der Weide K, de
27 Bont ES, Huls G, Vellenga E *et al*: **Gene expression profiling in the leukemic stem cell-enriched**
28 **CD34+ fraction identifies target genes that predict prognosis in normal karyotype AML.**
29 *Leukemia : official journal of the Leukemia Society of America, Leukemia Research Fund, UK* 2011,
30 **25**(12):1825-1833.
- 31 12. Hu G, Schones DE, Cui K, Ybarra R, Northrup D, Tang Q, Gattinoni L, Restifo NP, Huang S, Zhao K:
32 **Regulation of nucleosome landscape and transcription factor targeting at tissue-specific**
33 **enhancers by BRG1.** *Genome Res* 2011, **21**(10):1650-1658.
- 34 13. Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B: **Mapping and quantifying mammalian**
35 **transcriptomes by RNA-Seq.** *Nat Methods* 2008, **5**(7):621-628.
- 36 14. Liberzon A, Subramanian A, Pinchback R, Thorvaldsdottir H, Tamayo P, Mesirov JP: **Molecular**
37 **signatures database (MSigDB) 3.0.** *Bioinformatics* 2011, **27**(12):1739-1740.
- 38 15. Lin D: **An Information-Theoretic Definition of Similarity.** In: *Proceedings of the Fifteenth*
39 *International Conference on Machine Learning.* Morgan Kaufmann Publishers Inc. San Francisco,
40 CA, USA 1998: 296-304.
- 41 16. Frohlich H, Speer N, Poustka A, Beissbarth T: **GOSim--an R-package for computation of**
42 **information theoretic GO similarities between terms and gene products.** *BMC Bioinformatics*
43 2007, **8**:166.
- 44 17. Yang X, Regan K, Huang Y, Zhang Q, Li J, Seiwert TY, Cohen EE, Xing HR, Lussier YA: **Single sample**
45 **expression-anchored mechanisms predict survival in head and neck cancer.** *PLoS Comput Biol*
46 2012, **8**(1):e1002350.

- 1 18. Harris SL, Levine AJ: **The p53 pathway: positive and negative feedback loops.** *Oncogene* 2005,
- 2 **24(17):2899-2908.**
- 3 19. Xiong Q, Ancona N, Hauser ER, Mukherjee S, Furey TS: **Integrating genetic and gene expression**
- 4 **evidence into genome-wide association analysis of gene sets.** *Genome Res* 2012, **22(2):386-397.**
- 5