

## Supplementary Method

### Discovery Cohort:

A total of 35 T-ALL patients based on morphology, cytochemistry and immunophenotypic features were studied for RNA-Seq. The age of the patients ranged from 1 to year 55 with a median age of years. There were 26 males (19 pediatric, 7 adults) and 9 females (7 pediatric, 2 adults). The mean hemoglobin (Hb), total leukocyte count (TLC) and platelet count (Plt) was 8.5 gm/dl, 166,300/ $\mu$ L and 54,500/ $\mu$ L, respectively. The immunophenotypic subtype of these patients is shown in Table 1 . Transcriptome data from a pool of total RNAs from 5 normal human thymuses were used as control. (this was kindly provided by Prof. Jan Cools, Leuven, Belgium)

**Table 1: Distribution of T-ALL patients in discovery cohort based on immunophenotype**

| Subtype of T-ALL |             | Number of patients |
|------------------|-------------|--------------------|
| Immature T-ALL   | ETP-ALL     | 5                  |
|                  | Non ETP-ALL | 8                  |
| Cortical T-ALL   |             | 17                 |
| Mature T-ALL     |             | 5                  |

### Morphological diagnosis

In this study, the morphological evaluation was performed by hematopathologist, in accordance with (Vardiman, Thiele et al. 2009, Taylor, Xiao et al. 2017). For each case, routine well-prepared

1 Jenner-Giemsa-stained smears were evaluated. Cytochemistry for myeloperoxidase (MPO) was  
2 done in each case.

### 3 **Flow cytometric analysis**

4 All monoclonal antibodies were obtained from Beckman Coulter [BC], Hialeah, FL. The  
5 antibodies were conjugated to fluorescein isothiocyanate (FITC), phycoerythrin (PE), PE-Texas  
6 Red (ECD), PE-cyanin 5 and PE-cyanin 7 (PE-Cy7). The antibodies :*CD3*, *CD45*, *CD2*, *CD7*,  
7 *CD4*, *CD5*, *CD8*, *CD1a*, *CD13*, *CD33*, *CD117*, *HLA-DR*, *CD34*, *CD65* and *CD11b* were used in  
8 flow cytometry. All procedure, sample acquisition and analysis protocol were used as described  
9 in publication from our group.(Chopra, Bakhshi et al. 2014).

10 The subclassification of patients was done based on EGIL classification.(Bene, Castoldi et al.  
11 1995). ETP-ALL immunophenotype was characterized by an absence (<5% positive  
12 lymphoblasts) of *CD1a*, *CD8* expression , weak *CD5* expression with less than 75% positive  
13 lymphoblast, and expression of one or more of the myeloid or stem cell markers (*CD117*, *CD34*,  
14 *HLA-DR*, *CD13*, *CD33*, *CD11B* or *CD65*) on at least 25% of lymphoblast (Coustan-Smith,  
15 Mullighan et al. 2009).

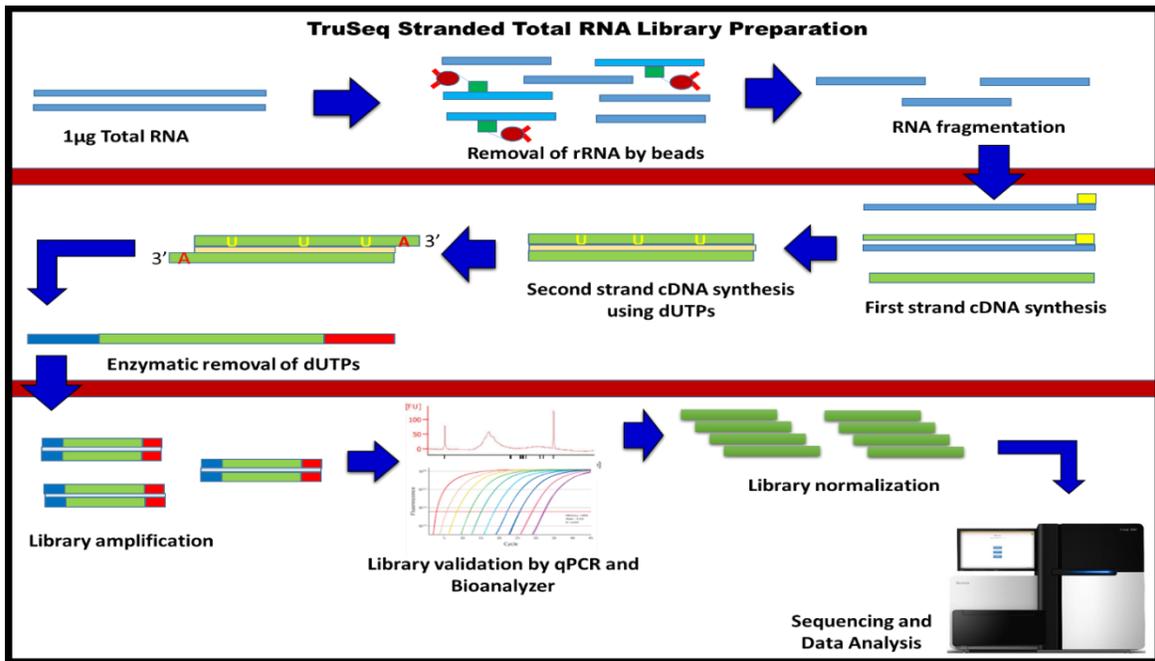
### 16 **RNA sequencing**

17 PBMC isolation performed by density gradient method using Histopaque (Sigma-Aldrich, USA)  
18 and 1 ml TRIzol (Thermo Fisher Scientific, USA) for per one million cells was added in each vial  
19 immediately and mixed by syringing. TRIzol mixed samples were stored at -80°C for RNA  
20 isolation. RNA was isolated by TRIzol (Guanidinium thiocyanate-phenol-chloroform extraction)  
21 method (Rio Dc Fau - Ares, Ares M Jr Fau - Hannon et al.). RNA was dissolved in DEPC-treated  
22 waterh a pipette tip. RNA quantity was checked by nanodrop and Qubit RNA broad range assay

1 kit on dye-based qubit spectrophotometer (Thermo Fisher Scientific, U.S.). Agilent 2100  
2 Bioanalyzer (Agilent Technologies, USA) with RNA chips was used to check the RNA Integrity  
3 Score (RIN). We selected samples with  $\geq 7$  RIN score for RNA sequencing library preparation.

#### 4 RNA seq library preparation

5 Sample preparation for sequencing was carried out using strand specific Truseq RNA sample  
6 preparation kit (Illumina, San Diego, California, U.S.) as per supplier's instructions and 8.0  
7 picomol of the pooled library was sequenced on the Illumina HiSeq2000 system (Figure 1 ).



8

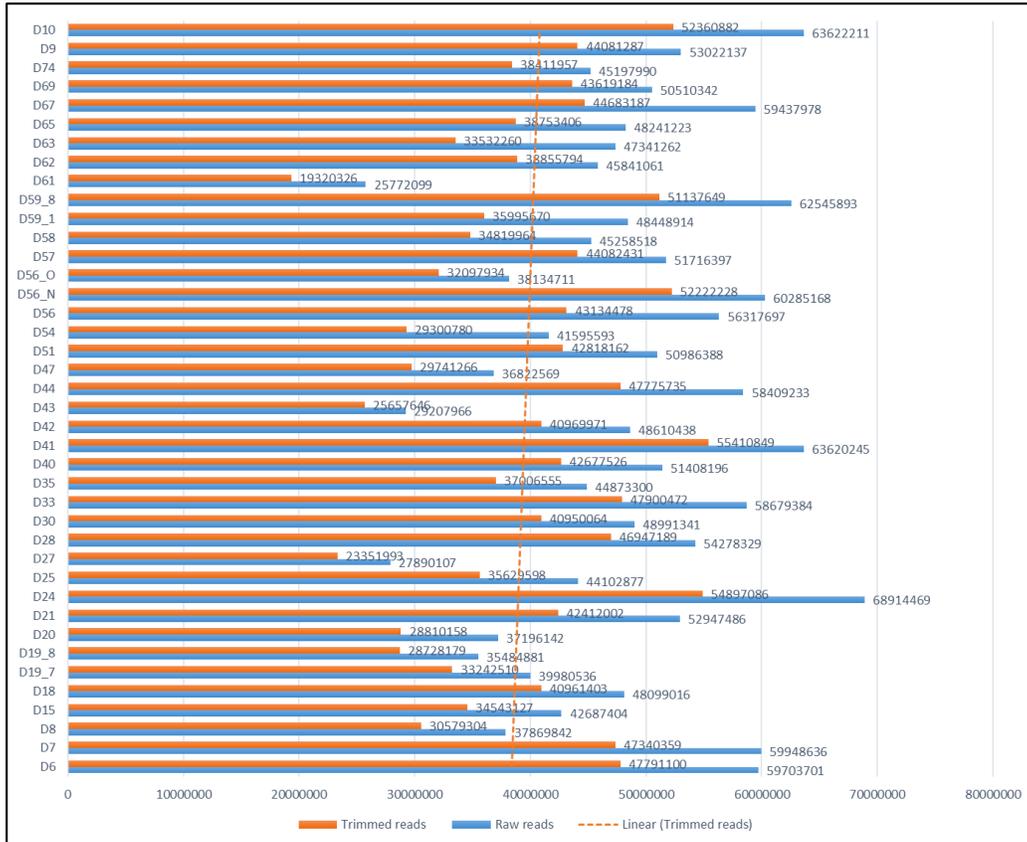
9 **Methods: Schematic diagram of library preparation by using Illumina TruSeq Stranded**

10 **Total RNA library preparation kit**

#### 11 Data collection and Quality analysis

12 Total of 480 GB raw data was collected in the form of “. fastq” file which contained 50 million  
13 average reads per file. The quality of reads was checked using FastQC v0.11.8 (Andrews 2010)

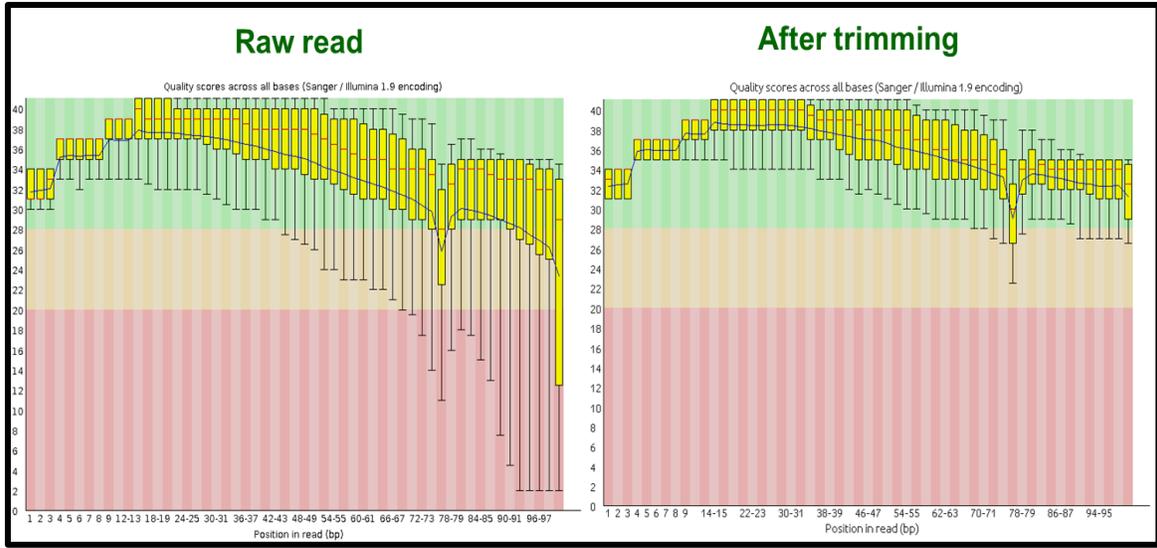
1 and trimming were implemented using Trimmomatic PE (Bolger, Lohse et al. 2014). We used a  
 2 quality filter of Phred quality cut-off 30 and ahead crop of 15. The overall percentage of reads  
 3 surviving after Trimmomatic was 70.83% to 86.36%.(Figure 2&3)



4

5

**Figure 2 : Graphical representation of raw and trimmed reads**



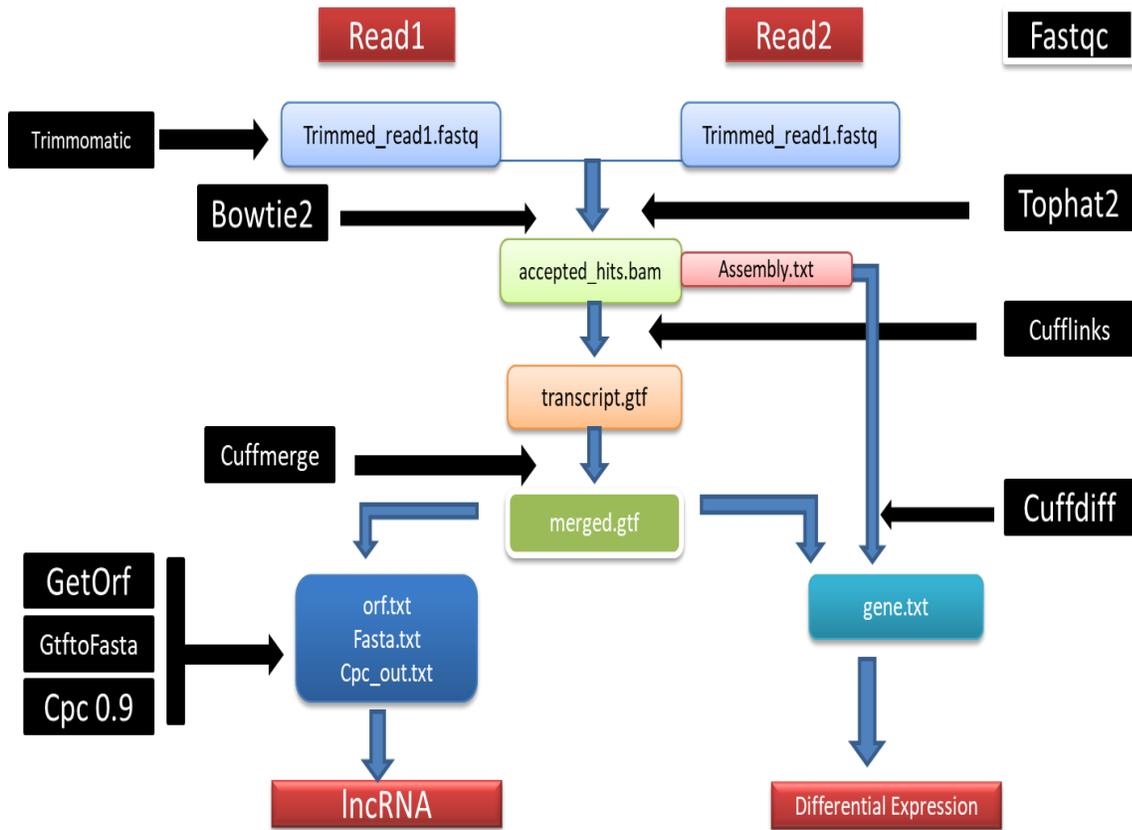
1

2

**Figure 3: Quality check of sample D6 raw reads and after running Trimmomatic software**

3

**by FastQC.**



4

1 **Figure 4: Data analysis pipeline for differentially expressed coding genes and long non-**  
2 **coding RNA**

3 **Mapping and alignment of reads**

4 Brief analysis pipeline was described in Figure 4. Human reference genome, hg38, was used for  
5 mapping and alignment, downloaded from UCSC genome browser gateway. Reads were aligned  
6 to the indexed hg38 genome using splice junction mapper: Tophat2.1.1 (Trapnell, Pachter et al.  
7 2009, Kim, Pertea et al. 2013). Total Avg 82.275 % reads were aligned to the human genome hg38.

8 **Transcriptome assembly and merging:**

9 Cufflinks 2.2.1 (Trapnell, Williams et al. 2010) was used to assemble the transcripts, estimate  
10 transcripts abundance and calculate differential gene expression and regulation. These assemblies  
11 .gtf files are then merged together into a single file by using the Cuffmerge utility. Finally, merged  
12 files were stored in merged.gtf format.

13 **Differential gene expression and functional analysis by the supervised approach**

14 Differential expression of the assembled transcripts across 3 types of T-ALL (immature, cortical  
15 and mature T-ALL based on immunophenotypic features) were analyzed using module Cuffdiff  
16 from the Cufflinks package (Trapnell, Williams et al. 2010). Samples of cortical subgroup were  
17 also analyzed in sCD3 positive and negative conditions by Cuffdiff. The Fragment per Kilobase  
18 Exon per Million Reads (FPKM) measure was used to quantify the expression of the transcript.  
19 For better understanding, these FPKM values were converted into fold changes. We filtered in all  
20 genes which were  $\geq 2$ -fold changes as a significantly important gene for a given condition.

1 Gene functional analysis was done on DAVID 6.8 database (Huang da, Sherman et al. 2009, Huang  
2 da, Sherman et al. 2009). Gene list of each subgroup was analyzed and result for “biological  
3 process, cellular component and KEGG pathways were filtered in, based on p-value  $\leq 0.005$ .

#### 4 **Differential gene expression analysis by the unsupervised approach**

5 The output of cufflink was used in cuffnorm to produce several output files that contain normalized  
6 fragment counts and expression levels at the level of transcripts, primary transcripts, and genes. In  
7 this approach, 35 patients and transcriptome from normal thymus was used. We found 58,224  
8 genes that were expressed in all 35 patient samples including normal thymus. In cuffnorm output,  
9 we removed all transcripts which were  $< 2$  FPKM scores and we found 31,694 genes that were  $\geq 2$   
10 fpkm. All samples were used as an input feature to perform principal component analysis (PCA)  
11 on BioVinci Version: 1.1.5, r20181005 (Bioturing, USA). We found maximum principle  
12 component equal to the samples and only 6 were statistically significant. To apply an unsupervised  
13 approach all 6 principal components were analyzed, and a graph was plotted. In this graph, all  
14 sample were clustered into 3 major clusters. After this, all samples were compared with  
15 immunophenotypic phenotypes of samples.

#### 16 **Differential expression analysis of non-coding RNA**

##### 17 **Known lncRNA**

18 The long noncoding RNAs were identified and characterized by “Biomart”  
19 (<https://www.ensembl.org/biomart>). Later their differential gene expression analysis was  
20 performed. All transcript which had more then 1FPKM were used to calculate fold changes.  
21 Transcripts those were more than 2-FC and less than -2 FC was selected as a putative transcript in

1 a given condition. The results were plotted on heatmap using multiple experiment viewer (MeV,  
2 <http://mev.tm4.org>). Further, their gene function co-expression network analysis was performed.

### 3 **Novel lncRNA**

#### 4 **Annotation of novel long non-coding RNA**

5 For identification of novel lncRNAs, the resulting output of cuffmerge was used to analyze  
6 lncRNA in the data. The transcripts with a nucleotide length greater than 200 were filtered. These  
7 transcripts were further classified as coding and non-coding by evaluating the length of their open  
8 reading frames (ORF) and coding potentials using two independent approaches, ORF finder and  
9 CPC calculator (Kong, Zhang et al. 2007). The ORFs of all the transcripts were predicted using  
10 getORF utility of EMBOSS toolkit. Transcripts having an ORF of length greater than or equal to  
11 30 amino acids (aa) were removed in order to remove all protein-coding as well as micro peptide  
12 coding transcripts, thus giving more stringent criteria to predict only noncoding transcripts as  
13 lncRNAs. The coding potential of the filtered transcripts was further calculated using the Coding  
14 Potential Calculator (CPC). A CPC scores less than zero indicates a low coding potential of a  
15 transcript while positive CPC score indicates a high coding potential of a transcript. The transcripts  
16 with a score less than 0 were retained and the final set represented the long non-coding RNAs of  
17 zebrafish while the one with orf more than 30 aa and CPC score of more than 1 were considered as  
18 protein-coding transcripts. The pipeline used in the analysis of novel lncRNA has been  
19 summarized in Figure 4.

20 Differential expression analysis of lncRNAs: Differential expression of the assembled transcripts  
21 across various types of T-ALL (immature, cortical and mature T-ALL) was analyzed using  
22 Cuffdiff. The FPKM measure was used to quantify the expression of lncRNAs. lncRNAs which

1 have at least >5 FPKM in one type of T-ALL and at least a 2-fold difference between other  
2 conditions were prioritized for further validation.

3

#### 4 **Differential gene expression of small non-coding RNA (miRNA)**

5 Like, lncRNAs, miRNAs were identified from “Biomart” (<https://www.ensembl.org/biomart>). To  
6 get the more information differential gene expression analysis was performed. All transcripts  
7 which were  $\geq$  to 1 fpkm were used to calculate fold changes. Transcripts with  $\geq$  2-FC were selected  
8 as putative miRNA transcript in given subgroup and heatmap was plotted using MeV. Due to  
9 insufficient sample volume, real-time validation could not be performed on RNA-Seq samples for  
10 miRNA.

#### 11 **Gene co-expression network analysis**

12 In order to predict the functions of the genes, histone modifiers, lncRNA, miRNA in all three  
13 subgroups a correlation network analysis was performed. A correlation matrix was generated using  
14 Pearson correlation as a statistical measure using R scripts based on the FPKM scores of the genes  
15 and package Hmisc where rcorr function was used. A threshold of “r” value 0.9 and “p-value” less  
16 than 0.05 was used to construct the co-expression network. Cytoscape (Shannon, Markiel et al.  
17 2003) was used to visualize the network where rscore was used as attribute measure. All DF genes,  
18 lncRNA, miRNA and histone modifiers were used as a source node and other genes as target  
19 source. The 3D-way layout was used for visualization of this network.

#### 20 **Data visualization**

1 All differentially expressed gene in all the subgroups were visualized by heatmap based on the log  
 2 scale and hierarchical clustering based on complete linkage as input settings. For gene network  
 3 analysis, Cytoscape tool (Shannon, Markiel et al. 2003) was used that allows to build a correlation  
 4 network in between the gene and other entity based on r and p-value. In network analysis, all gene  
 5 nodes were plotted using a 3D layout.

## 6 Validation of RNA seq results on the discovery cohort

7 Validation of data was done by reverse transcriptase PCR followed by real-time PCR using  
 8 specific primers and probes listed in Table 2 & 3 Primer3 (Soulie, Clappier et al. 2005).

9 cDNA synthesis was done using Vilo reverse transcriptase kit (Thermo Fisher Scientific, USA)  
 10 with 100ng of RNA. Further cDNA was diluted in 1:10 ration in nuclease-free water.

11

12

13 **Table 2: Primer sequence used for validation by real-time PCR.**

|                  |                         | Forward primer (5'-3')      | Reverse (5'-3')            |
|------------------|-------------------------|-----------------------------|----------------------------|
| <b>Immature</b>  | <i>MEF2C</i>            | GCGCTGATCATCTTCAAC          | CTTGCCTGCTGCTGATCATT       |
|                  | <i>BALLC</i>            | GCCCTCTGACCCAGAAACAG        | CTTTGCAGGCATTCTCTTAGCA     |
|                  | <i>LYL1</i>             | TGGCCCTGCACTACCACC          | AGATGCTAAAAGGTCCTGCTGG     |
|                  | <i>LMO2</i>             | CTGCCTGAGCTGCGACCT          | GTTGTAGTAGAGGGCCCGC        |
|                  | <i>HHEX</i>             | TCTGCATAAAAGGAAAGG          | CGTCTCGAATTTCTTCTC         |
|                  | <i>TRGC1</i>            | TACCTCCTCCTGCTCCTCAA        | GAGAACTGAAATGGCCCAA        |
| House<br>keeping | <i>ABL1</i>             | TGGAGATAAACTCTAAGCATAAAAGGT | GATGTAGTTGCTTGGGACCCA      |
|                  | <i>GAPDH</i>            | GAAGGTGAAGGTCGGAGTCAAC      | CAGAGTAAAAGCAGCCCTGGT      |
|                  | <i>β-actin</i>          | TCACCCACACTGTGCCCATCTACGA   | CAGCGGAACCGCTCATTTGCCAATGG |
|                  | <i>HOX11<br/>(TLX1)</i> | TGGAGAGTAACCGCAGATACACA     | CGTGCGGGCTTCTTC            |
|                  | <i>TRGV2</i>            | CCTCAGTCACTTCCCTCTGC        | GACTGGCAGGAGACAGGAAA       |

|            |                  |                        |                         |
|------------|------------------|------------------------|-------------------------|
| Cortical   | <i>TLX3</i>      | ACGGTCTCCAGCCTTGGC     | GCGCCGGGTCACGGT         |
|            | <i>FAT1</i>      | TGCTACAGACGCAGACATCC   | AACAGCTTGCTCCTCACGAT    |
|            | <i>MEG11</i>     | GGAGCTGAATCCCTACACCA   | GTCTCGGCCCTTCTCTTTCT    |
|            | <i>LMO1</i>      | CAACGTGTATCACCTCGACTGC | GAAGAATTTGTCTCCACACAAAA |
|            | <i>FGR</i>       | GCATCCAAGGTGTGGAGTTT   | GCAACGAAGGGGACATTTTA    |
|            | <i>TTK</i>       | CAGCAGCAACAGCATCAAAT   | TGCTTGAACCTCCACTTCT     |
| Mature     | <i>ST20</i>      | CTCTGTCTCCCAGGTTTCAGG  | TGCTTCAGGACAGAAGAGCA    |
|            | <i>EML4</i>      | CCATGCAACGAGATAAGCAA   | CACATGCAGCTGAAGGAAAA    |
|            | <i>CDCA2</i>     | GCCCTGCACTGTATCGAAAT   | CGCTGAGACCTTCTTTTCTG    |
|            | <i>TAL1</i>      | TTCACCACCAACAATCGAGTG  | TGCCCCATCGCTCC          |
|            | <i>CREG2</i>     | TCTCTGCCTCCTGATCCTGT   | GGCAAAGTCTTTGAAGCAG     |
| lncRNA     |                  |                        |                         |
|            | <i>XIST</i>      | AAGCTAAGGGCGTGTTCAGA   | TGACTTCTCTGCCTGACCT     |
|            | <i>LINC01221</i> | TTCCAGGCAAATCCTTTTGT   | ACCAGAGAAGCCTGCCAGTA    |
|            | <i>TRBV11-2</i>  | CTCCTGGGAGCAGAACTCAC   | TAAAGGGTAGCATGGCCAGA    |
|            | <i>MALAT1</i>    | TGTGTGCCAATGTTTCGTTT   | AGGAGAAAGTGCCATGGTTG    |
|            | <i>CYTOR</i>     | CGGAATGCAGCTGAAAGATT   | ACCAGCCCATGACCAAAATA    |
|            | <i>HOTAIR</i>    | AAGGCCCAAAGAGTCTGAT    | AGCACTTCTCTCGCCAATGT    |
|            | <i>HOTTIP</i>    | AAGGGTCTCAGCTCCACAGA   | CTGCCGTCTTTTCTGAGTCC    |
| Epigenetic | <i>RAG1</i>      | AGCCTGCTGAGCAAGGTACC   | GAAGTGAAGTCCCAAGGTGGG   |
|            | <i>EZH2</i>      | TTCATGCAACACCCAACACT   | GAGAGCAGCAGCAAACCTCT    |
|            | <i>KDM6A</i>     | CCTCATAACCGCACAAACCT   | CTGCCGAATGTGAACTCTGA    |
|            | <i>EP300</i>     | CAAACGCCGAGTCTCTTTTC   | GTTGAGCTGCTGTTGGCATA    |
|            | <i>SETD2</i>     | TCACAAGGCAGACTCAGTGG   | CTGCTGTCTTGGGCTTTTTC    |

1 **Table 3: Probe sequence for real time validation**

| <b>GENE</b>        | <b>5'-PROBE-3'</b>                       |
|--------------------|--|
| <i>MEF2C</i>       | FAM- CAAGCTGTTCCAGTATGCCA-BHQ1           |
| <i>LYL1</i>        | FAM-TCACCCCTTCTCAACAGTGTCTACATTGG-BHQ1   |
| <i>HHEX</i>        | FAM-AGATTCTCCAACGACCAGACCATC-BHQ1        |
| <i>TRGC1</i>       | TET-AGAAGAACGGCTTTCTGCTG-BHQ1            |
| <i>β-actin</i>     | QUASAR705-ATGCCCTCCCCATGCCATCCTGCGT-BHQ1 |
| <i>HOX11(TLX1)</i> | TET-AGGACAGGTTACAGGTCACCCCTATCAGAA-BHQ1  |
| <i>TRGV2</i>       | CY5-TGCTCCTCTTCATCTGGTCC-BHQ1            |
| <i>FAT1</i>        | TxRED-ACGTTATTGGGTTTCAGGTGC-BHQ1         |
| <i>MEG11</i>       | TET-CAGGCATCATGCTCCTGTTA-BHQ1            |
| <i>FGR</i>         | CY5-CTGCTCATCTGCACGTGAGT-BHQ1            |
| <i>TTK</i>         | TET-CAGCAAATGAATGCATTTTCG-BHQ1           |
| <i>ST20</i>        | TxRED-GCCTAAATCTGGCTGGATGA-BHQ1          |
| <i>EML4</i>        | TxRED-GTCCTAACCCCTGGCTTCA-BHQ1           |
| <i>CDCA2</i>       | FAM-CTGGAATTCTCAGAGGCAGG-BHQ1            |

|                  |  |
|------------------|--|
| <i>TAL1</i>      | TET-ATTACTGATGGTCCCCACACCAAAGTTGTGC-BHQ1 |
| <i>CREG2</i>     | CY5-TAAGGTCCCCTGGCTAGGTT-BHQ1            |
| <i>XIST</i>      | TET-CCATTTTAGACTTGCCGCAT-BHQ1            |
| <i>LINC01221</i> | TXRED-CAGTCTGAACAGCCAGTGA-BHQ1           |
| <i>TRBV11-2</i>  | CY5-GAGAAAAGGCAGAGTGTGGC-BHQ1            |
| <i>MALAT1</i>    | QUASAR705-CAGCAGGCAGCTGTAAACAG-BHQ1      |
| <i>CYTOR</i>     | TXRED-CTGTGGACTCTGAGGCCTCT-BHQ1          |
| <i>HOTAIR</i>    | CY5-AAAGGCTGAAATGGAGGACC-BHQ1            |
| <i>HOTTIP</i>    | CY5-GAGGAGCCACAAGCAGGTAC-BHQ1            |
| <i>EZH2</i>      | CY5-TTACCAGCATTGGAGGGAG-BHQ1             |
| <i>KDM6A</i>     | TXRED-AGCCGTGGAAAAACCAACTA-BHQ1          |
| <i>EP300</i>     | TET-TCAATGTTGCATTCAGCCAT-BHQ1            |
| <i>SETD2</i>     | QUASAR705-TGAGATGGACCTGGGAACTC-BHQ1      |

1

2 **Validation cohort:**

3 **Patients characteristics**

4 A total of 99 T-ALL patients including 70 children and 29 adults were studied for validation of  
5 RNA-Seq results. The age of the patients ranged from 3 to 65 years with a median age of 12 years.  
6 There were 86 males (61 pediatric, 25 adults) and 13 females (9 pediatric, 4 adults). The mean  
7 hemoglobin (Hb), total leukocyte count (TLC) and platelet count (Plt) was 9.5 gm/dl, 104,922/ $\mu$ L  
8 and 87,895/ $\mu$ L, respectively. The immunophenotypic subset of these patients are shown in Table  
9 1 of the 99 T-ALL patients in the validation cohort, RNA in enough quality and quantity was  
10 available for the determination of MEF2C in 99 and BAALC, HHEX, LYL1, TAL1, FAT1, XIST  
11 genes in 87 cases.

12 **Table 3: Distribution of T-ALL patients in validation of based on immunophenotype into**  
13 **immature, cortical and mature subtypes.**

| Subtype of T-ALL | Number of patients |
|------------------|--------------------|
|                  |                    |

|                |             |    |
|----------------|-------------|----|
| Immature T-ALL | ETP-ALL     | 13 |
|                | Non ETP-ALL | 32 |
| Cortical T-ALL |             | 44 |
| Mature T-ALL   |             | 10 |

1

2 The results of RNAseq data were validated in the validation cohort using realtime PCR to assess  
3 its clinical and prognostic significance.

4 **Treatment**

5 In the discovery cohort, 14 patients were treated with ICICLE protocol, 3 with Berlin-Frankfurt-  
6 Munster-90 (BFM-90) protocol, 1 with INCTR, 1 with Holzer’s protocol and 3 with hyper CVAD.  
7 13 patients did not take treatment.

8 In the validation cohort, seventy-two patients were treated with ICICLE protocol, 15 with Berlin-  
9 Frankfurt-Munster-90 (BFM-90) protocol and 3 with hyper CVAD. Nine patients did not take any  
10 treatment. Two patients died during induction chemotherapy. Complete remission was defined as  
11 bone marrow blasts <5% with a recovery of blood counts at the end of 4 weeks of induction  
12 chemotherapy. Any failure to do so (including the persistence of leukemic blasts in an  
13 extramedullary site), or death during induction therapy due to any cause, was considered as  
14 induction failure Patients who failed with one protocol were re-induced with another.

15 **Statistical analysis**

1 Fisher's exact test for categorical data was used to compare baseline clinical variables across  
2 groups in the validation cohort. A P value  $\leq 0.05$  (two-sided) was considered significant. *BAALC*,  
3 *LYLI*, *FAT*, *TALI* and *HHEX* were dichotomized at its median with patients classified as low gene  
4 expression if they had expression values within the lower 50% and as high gene expression if they  
5 had gene expression values within the upper 50%. For *MEF2C* gene expression, the patients were  
6 dichotomized into low *MEF2C* expression and high *MEF2C* expression based on a cutoff of 0.1-  
7 *MEF2C* high, if the expression is  $\geq 0.1$  and *MEF2C* low if the expression is  $< 0.1$ . This cut-off was  
8 chosen based on our unpublished data from a larger number of patients, which showed that the  
9 patients can be divided into groups with clinically distinct event-free survival and overall survival.  
10 For *XIST* gene expression, we took the cut-off of 1. Since the *XIST* gene is present on X-  
11 chromosome, we did analysis for all the patients and only male patients separately, to see the  
12 association of various parameters and *XIST* gene expression.

13 Response to prednisolone was assessed by the examination of peripheral blood at day 8 of initiation  
14 of prednisolone therapy. The patients were divided into prednisolone sensitive and resistant based  
15 on the presence of  $< 1000/\mu\text{L}$  and  $\geq 1000/\mu\text{L}$  blasts, respectively.

16 Achievement of complete remission (CR) was assessed after completion of induction  
17 chemotherapy and presence of absolute neutrophil count  $> 1000/\mu\text{L}$ , platelet count  $> 100,000/\mu\text{L}$   
18 and hemoglobin  $> 10\text{gm/dL}$ , no peripheral blood blasts, less than 5% BM blasts and absence of  
19 extramedullary leukemia. Relapse was defined as the reappearance of blasts in the peripheral  
20 blood, more than 5% BM blasts, or development of extramedullary leukemia.

1 The median follow-up was 16 months. After excluding 13 patients in validation cohort who did  
2 not receive treatment, 90 patients were included for the analysis of event-free survival and overall  
3 survival.

4 Event-free survival (EFS) was defined as the time from diagnosis to the date of the last follow-up  
5 in complete remission or the first event (i.e., induction failure, relapse, secondary neoplasm, or  
6 death from any cause). Failure to achieve remission due to non-response was considered an event  
7 at time zero. Survival was defined as the time from diagnosis to death or the last follow-up. Patients  
8 lost to follow-up were censored at the last contact. The last follow up was carried out on April  
9 2019. The Kaplan-Meier method was used to estimate survival rates, with the differences  
10 compared using a two-sided log-rank test. Cox proportional hazard models were constructed for  
11 EFS and OS and used for univariate and multivariate analyses. Covariates included in the full  
12 model of OS and EFS were Sex, WBC ( $<100 \times 10^9/L$ ,  $>100 \times 10^9/L$ ), and age ( $<10$  years vs.  $>10$   
13 years), gene expression, immunophenotype, CNS involvement, response to prednisolone  
14 treatment, bone marrow remission status and presence of minimal residual disease after the end of  
15 induction chemotherapy. All analyses were performed using the SPSS statistical software package,  
16 version 20.0/STATA software, version 11.