# Supplemental Information

# Tracing the Derivation of Embryonic

# Stem Cells from the Inner Cell Mass

# by Single-Cell RNA-Seq Analysis

Fuchou Tang, Catalin Barbacioru, Siqin Bao, Caroline Lee, Ellen Nordman, Xiaohui Wang, Kaiqin Lao, and M. Azim Surani

**Supplemental Inventory**

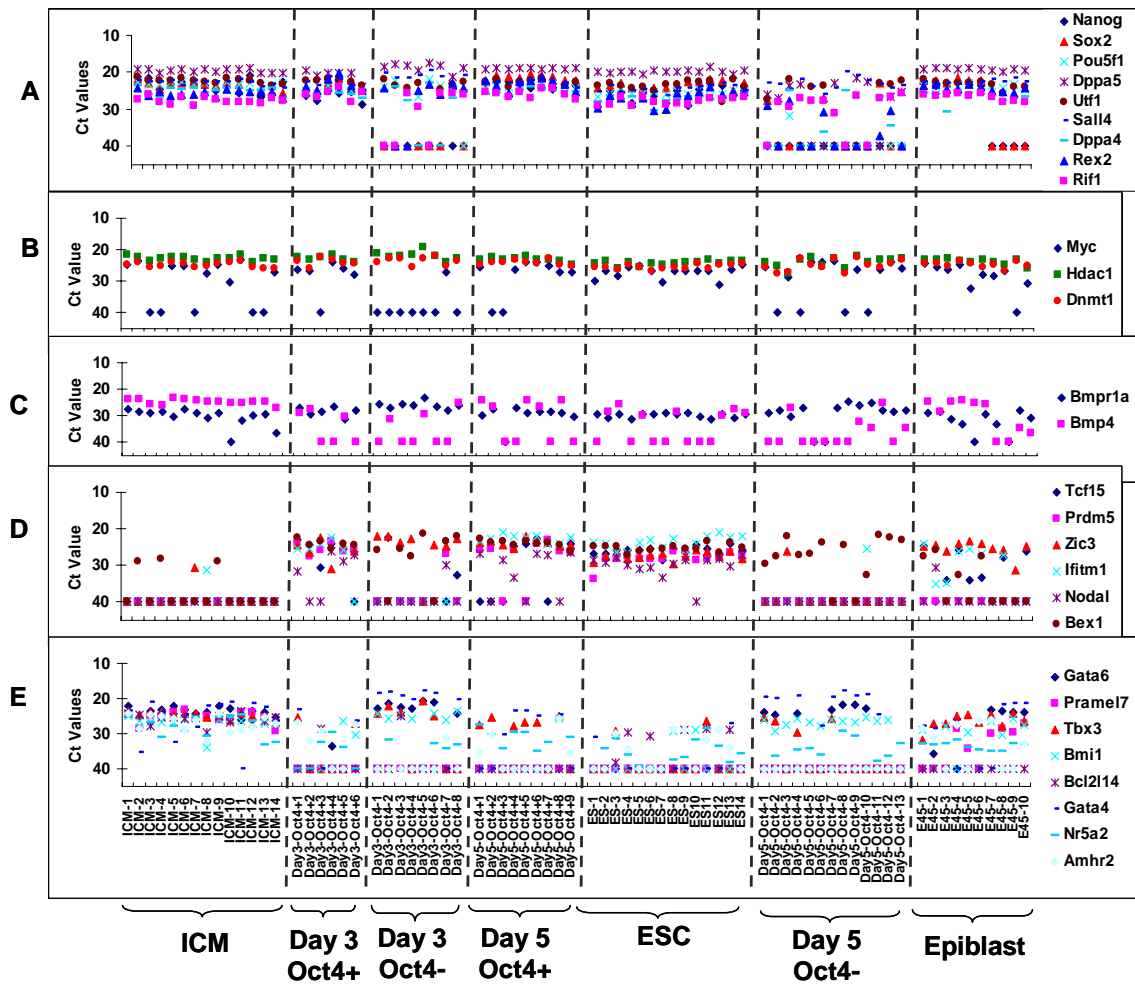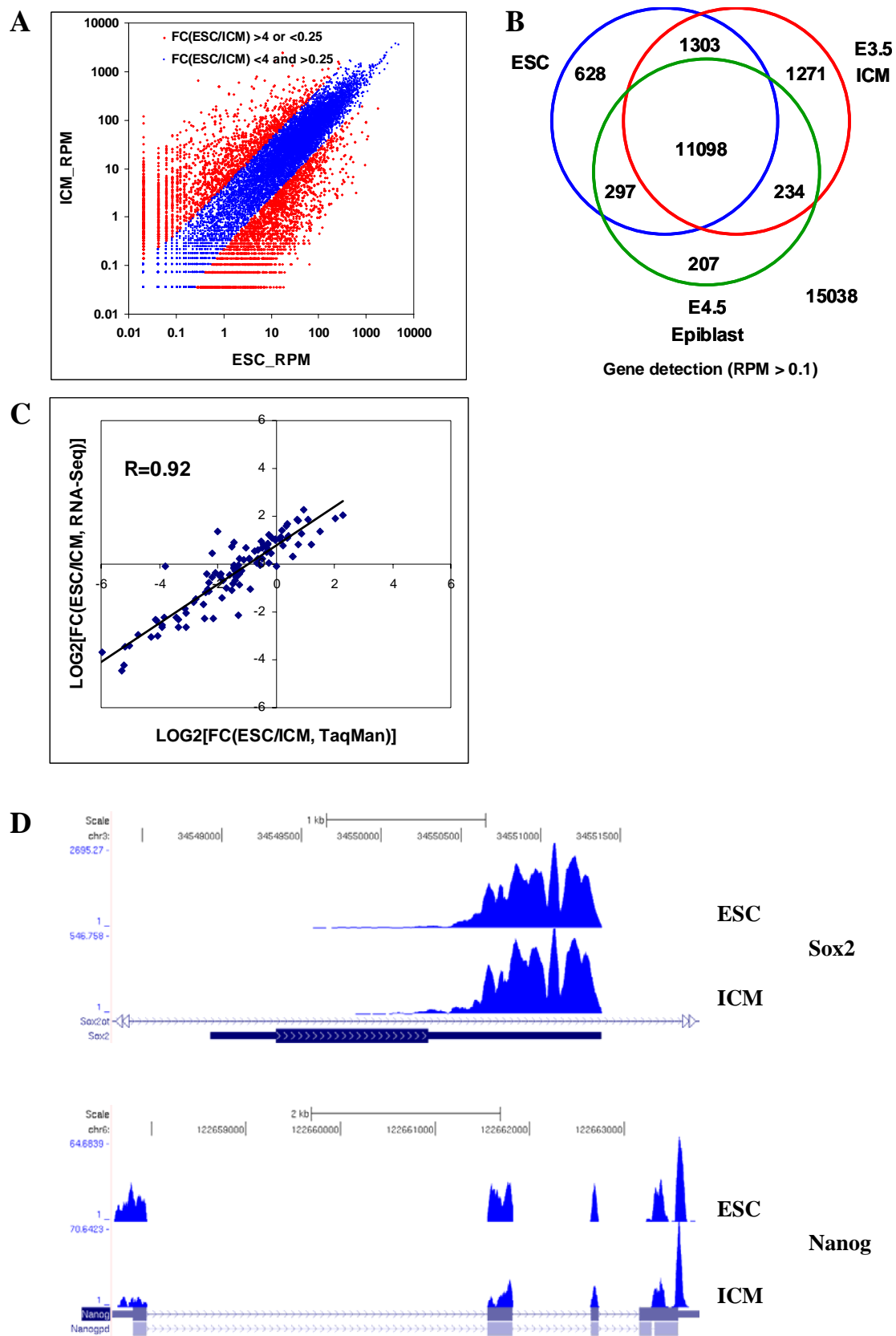| Items | Related to | General description |
|---|---|---|
| Figure S1 | Figure 2 | Gene expression dynamics measured by real-time PCR in single ICM outgrowth cells. The heterogeneity of ICM outgrowth, and the intermediate states of cells when they lose pluripotency were shown at single cell resolution. |
| Figure S2 | Figure 3 | Correlation plots and Pearson correlation coefficients of RNA-Seq data of single ICM outgrowth cells. The gene expression dynamics during ICM outgrowth at whole genome scale were shown. |
| Figure S3 | Figure 3 | Gene expression pattern between ICM, epiblast, and ESCs. The general overlap of the genes expressed between these types of cells was shown. |
| Figure S4 | Figure 5 | Splice specific differential expression of Dppa4 in ICM, epiblast, and ESCs. The dramatic regulation of alternative splicing was shown at single base resolution. |
| Figure S5 | Figure 6 | Gene network analysis of different pathways. Several of the gene networks were validated within an individual cell. |
| Figure S6 | Figure 3 | Gene expression dynamics during ICM outgrowth. The expression pattern of the genes related to pluripotency and self-renewal ability was shown. |
| Figure S7 | Figure 3 | Quality check of the single cell RNA-Seq data for ICM outgrowth. These analyses showed that the single cell RNA-Seq data are highly reproducible, reliable, and accurate for ICM outgrowth and ESCs. |
| Table S1 | Figure 1 & Figure 2 | The expression of 385 ESC related genes during ICM outgrowth by single cell real-time PCR. |
| Table S2 | Figure 3, Figure 4, Figure 5, & Figure 6 | Single cell RNA-Seq RefSeq transcripts and splice variants counts for ICM outgrowth. |
| Table S3 | Figure 3 | Summary of RNA-Seq alignments for ICM outgrowth. |
| Table S4 | Figure 6 | GO analysis of functions, networks, and pathways for ICM outgrowth. |
| Table S5 | Figure 7 | MicroRNA expression profile of ICM from E3.5 blastocysts and ESCs. |
| Table S6 | Figure 7 | MicroRNA enrichment calculations using target prediction algorithms of PicTar, Miranda, or TargetScan. |
| Table S7 | Figure 3 & Figure 7 | The lists of 1,378 potential pluripotency genes. |

**Figure S1**

**Figure S1** Gene expression dynamics measured by Real-time PCR in single cells of fourteen ICM cells (E3.5), ten epiblast cells (E4.5), six Day3 Oct4 positive ICM outgrowth cells, eight Day3 Oct4 negative ICM outgrowth cells, nine Day5 Oct4 positive ICM outgrowth cells, thirteen Day5 Oct4 negative ICM outgrowth cells, and fourteen ESCs: (A) *Oct4, Sox2, Nanog, Dppa4, Dppa5, Sall4, Utf1, Rex2*, and *Rif1* (B) *Myc, Dnmt1*, and *Hdac1*; (C) *Bmp4* and *Bmpr1a*; (D) *Tcf15, Prdm5, Zic3, Ifitm1, Nodal*, and *Bex1*; (E) *Gata4, Gata6, Pramel7, Tbx3, Bmi1, Bcl2l14, Nr5a2*, and *Amhr2* (Table S1).

**A**
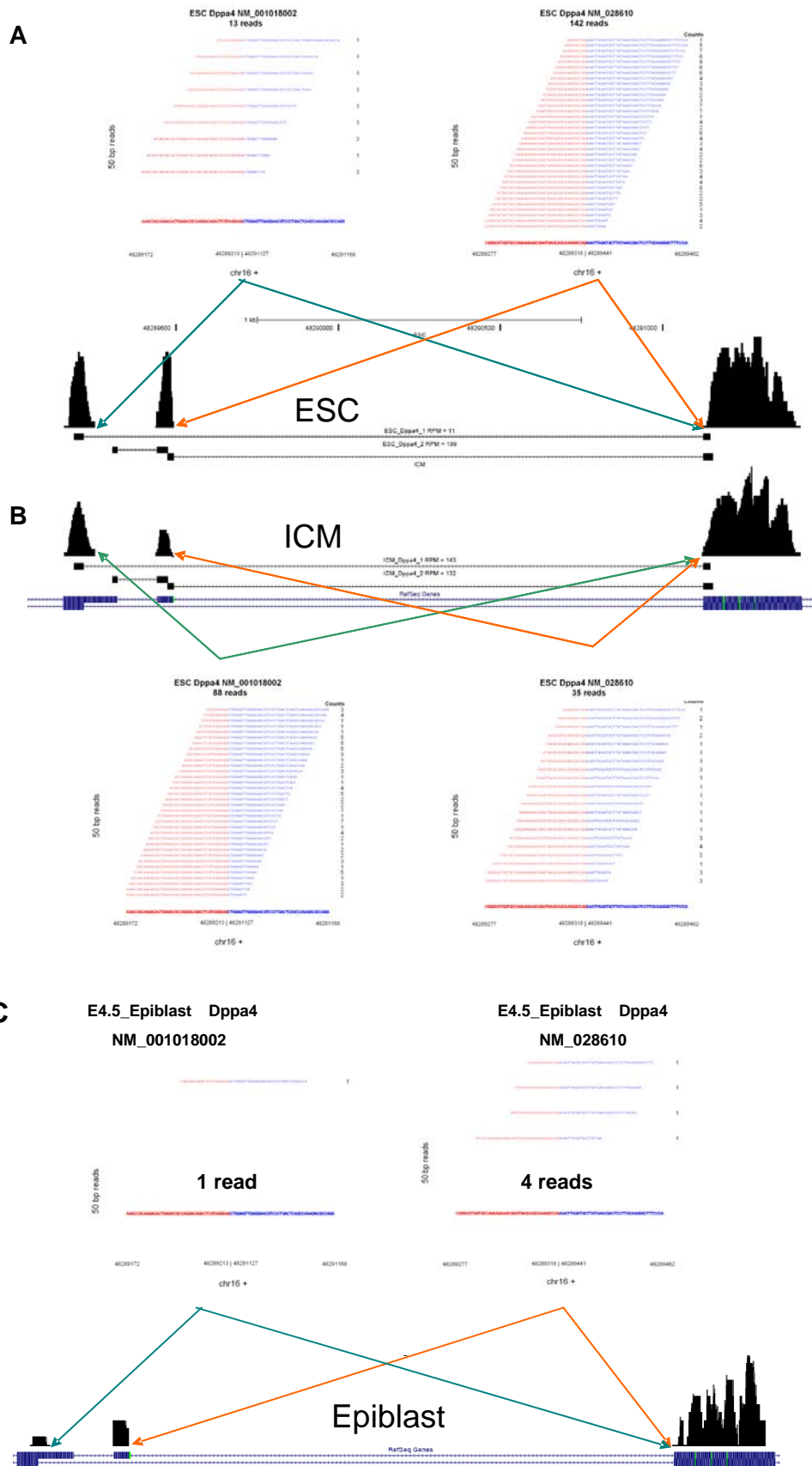


**B**



**C**



**Figure S2**

**Figure S2** Correlation plots and Pearson correlation coefficients of (A) 12 single ESCs, (B) 9 single ICM cells, and (C) 12 representative ICM outgrowth cells (Table S2).

**A**



**B**



Gene detection (RPM > 0.1)

**C**



**D**



**Figure S3**

**Figure S3** Gene expression pattern between ICM, epiblast, and ESCs. (A) Correlation plots for ESC vs ICM. The blue and red dots are, respectively, genes with fold changes less than or greater than 4; (B) Venn diagram of the genes expressed in ESC, E3.5 ICM, and E4.5 Epiblast (Table S2); (C) The correlation plot of the fold changes that are determined by RNA-Seq, LOG2 FC[ESC/ICM], and TaqMan real-time PCR, LOG2 of FC[ESC/ICM]; (D) Cover plots of Sox2 and Nanog in ESCs and ICM.

**A**

ESC Dppa4 NM_001018002
13 reads

ESC Dppa4 NM_028610
142 reads

ESC

**B**

ICM

ESC Dppa4 NM_001018002
88 reads

ESC Dppa4 NM_028610
35 reads

**C**

E4.5_Epiblast  Dppa4
NM_001018002

E4.5_Epiblast  Dppa4
NM_028610

1 read

4 reads

Epiblast

**Figure S4**

**Figure S4** Splice specific differential expression of Dppa4 in ICM, epiblast, and ESCs. RNA-Seq reads aligned to regions (exonic or exon-exon junctions) specific to a single transcript are used to differentiate between spliced transcripts. Dppa4 gene is known to encode two alternative spliced forms NM_001018002 (transcript variant #1, which has one specific exon and 2 specific exon-exon junctions, one of them being represented in this plot) and NM_028610 (transcript variant #2, which has one specific exon-exon junction). The number of reads covering NM_001018002 specific junction in (A) ESC (13 counts) is much less than the coverage obtained from (B) ICM cells (88 counts). The opposite trend is observed for NM_028610 specific exon-exon junction, for which we observe more reads in ESC (142 counts) compared to the reads observed the ICM (35 counts, Table S2).
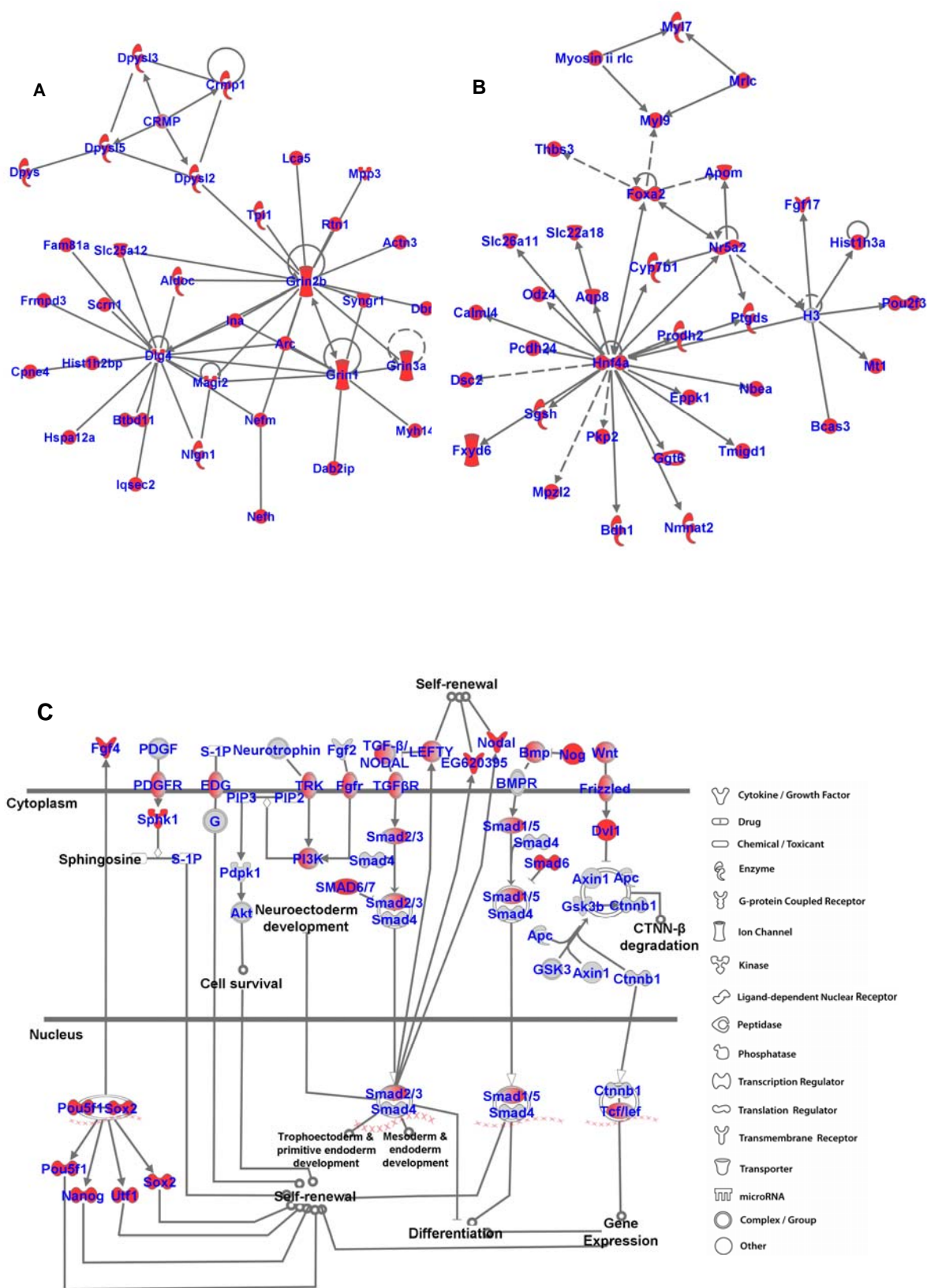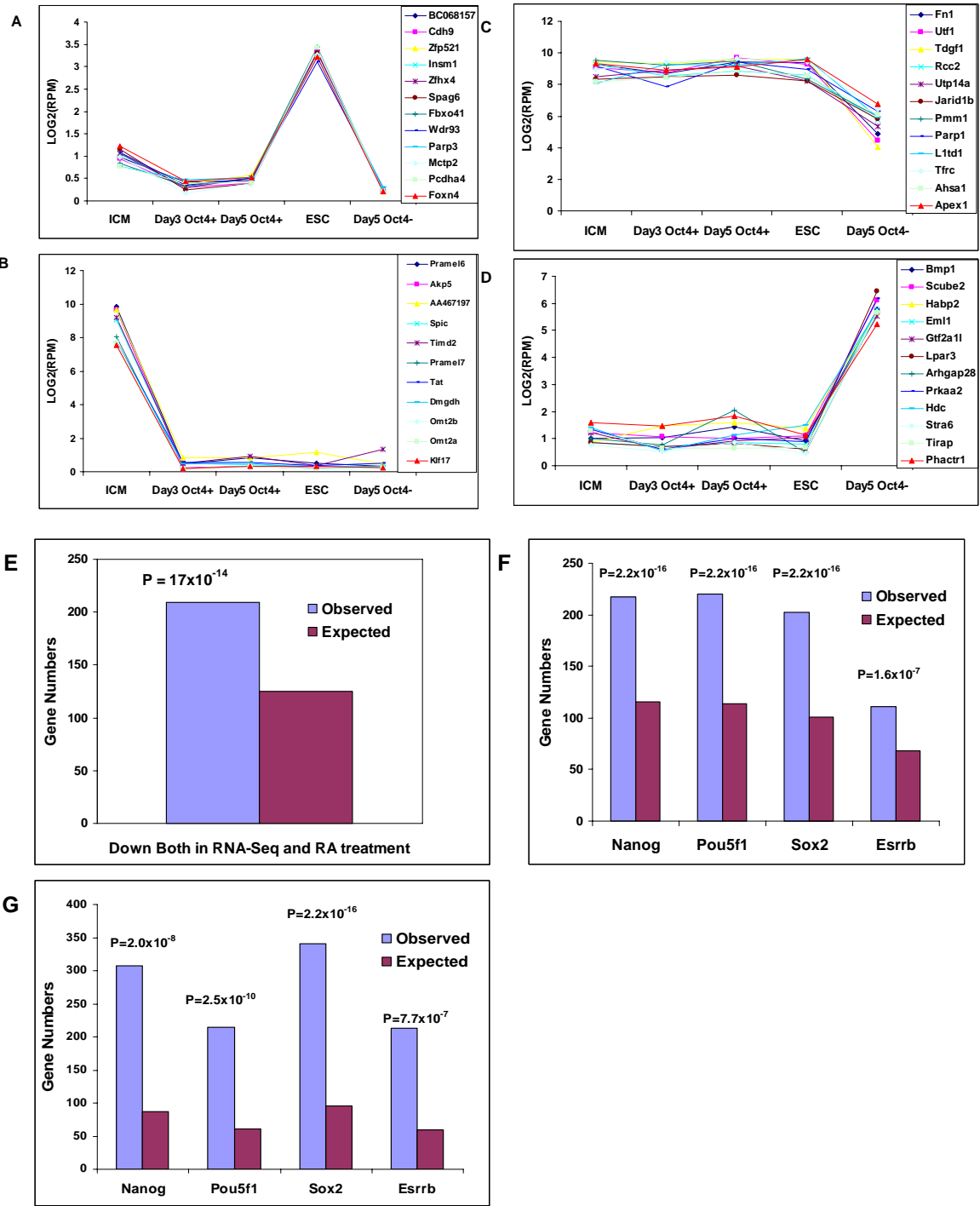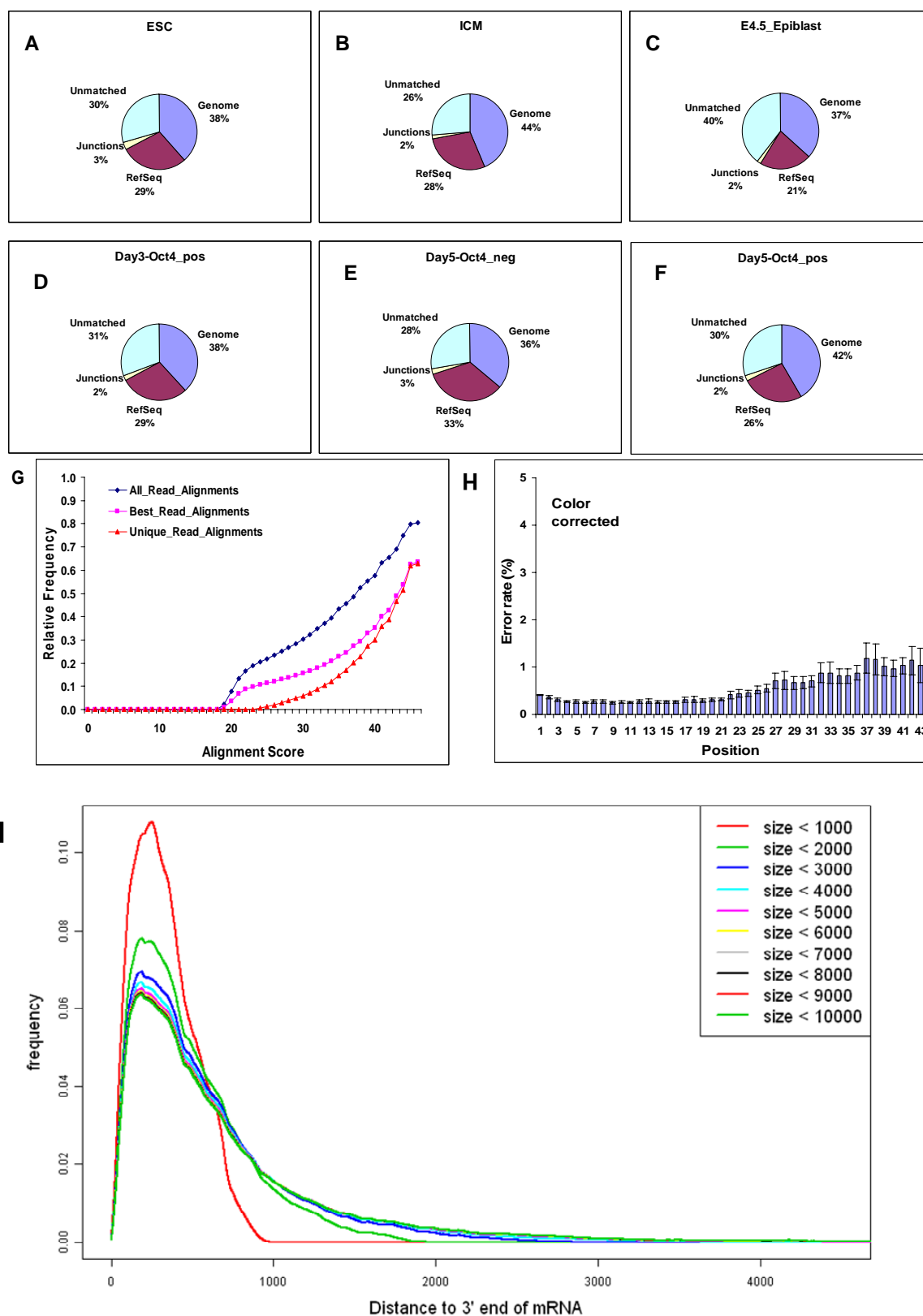
**Figure S5**

**Figure S5** Gene network analysis of different pathways. (A) Gene network plots of nucleic acid metabolism genes; (B) Gene network plots of lipid metabolism genes, which show dynamic changes of expression between ICM and ESC (FC[ESC/ICM]>4 or <0.25, p<0.01, Table S4). (C) Gene network plots of stem cell pluripotency pathway. 45 out of the 137 genes in this network are up/down regulated for more than 4 folds when the pluripotent cells lose pluripotency (The genes of FC[Day5-Oct4$^+$/Day5-Oct4$^-$]>4 or FC[Day5-Oct4$^+$/Day5-Oct4$^-$]<0.25, p<0.01 are shown in red color, Table S4). The p-value was estimated using Ingunity systems software (www.ingenuity.com).

**Figure S6**

**Figure S6** Gene expression dynamics during ICM outgrowth. (A) - (D), Candidate genes related to pluripotency and self-renewal ability of ESCs. (A) Expression profiling of the representative 12 genes that have low expression in E3.5 ICM, Day3 Oct4$^+$, and Day5 Oct4$^+$ ICM outgrowth cells, but are upregulated significantly in ESCs; (B) Expression profiling of the representative 11 genes that have high expression in E3.5 ICM, then decrease significantly at Day3 Oct4$^+$, Day5 Oct4$^+$ and Day5 Oct4$^-$ ICM outgrowth cells. Here, the averaged intensities of three single cells were plotted. (C) Expression profiling of the representative 12 genes that have high expression in ICM, Day3 Oct4$^+$ and Day5 Oct4$^+$ ICM outgrowth cells, and ESCs, but are downregulated significantly in Day5 Oct4$^-$ ICM outgrowth cells with no pluripotency; (D) Expression profiling of the representative 12 genes that have low expression in ICM, Day3 Oct4$^+$ and Day5 Oct4$^+$ ICM outgrowth cells, and ESCs, but are upregulated by more than 4 folds in Day5 Oct4$^-$ ICM outgrowth cells with no pluripotency. The averaged intensities of three single cells were plotted here (Table S2). (E) - (F): Enrichment of downstream genes of pluripotency master genes. (E) The earliest downregulated genes when the differentiation of ESC is induced by all-trans retinoic acid (RA) are overrepresented in the downregulated genes when pluripotent cells lose pluripotency in our assay. (F) The functional downstream genes of *Nanog, Oct4, Sox2*, and *Esrrb* that are down regulated both after the knockdown of these genes and when pluripotent cells lose pluripotency. (G) The functional downstream genes of *Nanog, Oct4, Sox2*, and *Esrrb* that are up regulated both after the knockdown of these genes and when pluripotent cells lose pluripotency. The number of the expected genes is equal to the number of the genes up- or down- regulated when one of the pluripotent gene is knocked down (FC>1.5), divided by the total number of the genes detected by RNA-Seq (RPM>0.1), and multiplied by the number of the genes up- or down- regulated in Day5 outgrowth when the pluripotent cells lose pluripotency (FC >2 or <0.5).

Figure S7

**Figure S7** Quality check of the single cell RNA-Seq data for ICM outgrowth. (A) - (F): Pie charts of the number of RNA-Seq reads. Reads are mapped to the mouse genome (mm9, NCBI Build 37) as described by Tang et al (2009) in the Alignment and Algorithm section. We used UCSC annotation database (mm9) to determine if matching locations of individual reads correspond to exon regions, or exon-exon junctions of known transcripts. The number of these reads as a fraction of the total number of reads produced by cell type is represented in these pie charts. We generated 500 millions reads in total (Table S3). We obtained about 60 - 74% of reads that mapped uniquely to Refseq, known junctions, and the genome. There are about 2 - 3 % reads that mapped uniquely to known exon junctions. Here, Genome means these reads uniquely mapped to the mouse genome, but did not map to RefSeq transcripts. The percentages are calculated based on aggregated reads for each cell type. The detailed alignment summary for each individual cell is listed in Table S3. (G) Alignment score for RNA-Seq reads. For each read alignment we calculate the score obtained by adding 1 for a color match and subtracting 1 for a mismatch. This scoring function is used in the extension step (described by Tang et al (2009) in the Alignment and Algorithm section). Cumulative distributions of all alignment scores, best read alignment, and alignment scores of reads aligning to a unique location, from all ESC cells, are shown. Low scoring alignments are expected to be produced by reads aligning over splice junctions or low quality reads, while high scoring reads are expected to represent exonic regions. Here, we only used unique reads for our RefSeq transcripts and junctions mapping. (H) The estimated sequencing error rate of SOLiD system. Reads that aligned contiguously to the genome (as described in Alignment and Algorithm section), on full length and to a unique location are used to estimate instrument error rate as a function of position within the read (see Tang et al Error detection section). For a given position, the number of times we see a difference between the read call and corresponding genomic location, calculated as a fraction of all reads considered, is represented on the y-axis. Data produced from all ESC and ICM cells is aggregated and represented as average error rates, and 95% confidence intervals of these estimates are inferred. The averaged error rate is about 0.5% for the first 30 positions. For 50-mer reads, we only plotted the first 44 bases since the extension step reaches full length of the read if

reads have fewer errors on the last positions (and so the error rates of the last positions for the subset of reads we used are under estimated). (I) Base coverage of the single cell RNA-Seq assay in single ES cells. To obtain the coverage length distributions of our cDNAs, we binned all 24,435 transcripts based on their sizes, with bin n containing all transcripts of size less than n kb. Base coverage is generated for each bin of transcripts and scaled to the total number of aligned reads. The obtained distribution is represented as a function of the distance to the 3' end of the transcripts. The read distribution for regions 3 kb away from the 3'end is very limited, which agrees with our gel results (data not show).

**Table S1** The expression of 385 ESC related genes during ICM outgrowth by single cell real-time PCR.

**Table S2** Single cell RNA-Seq counts (1) mapped to RefSeq for ICM, ESC and ICM outgrowth; (2) mapped to unique known exon-exon junctions and exons for each known splicing variant in RefSeq.

**Table S3** Summary of RNA-Seq alignments for ICM, ESC and ICM outgrowth.

**Table S4** GO analysis results of functions, networks, and pathways for (1) Refseq transcripts with FC[ESC/ICM]>4 or FC[ESC/ICM]<0.25, p<0.01); (2) Refseq transcripts with FC[Day5-Oct4$^+$/Day5-Oct4$^-$]>4 or FC[Day5-Oct4$^+$/Day5-Oct4$^-$]<0.25, p<0.01); (3) splicing variances with FC[ESC/ICM, junctions]>2 or FC[ESC/ICM, junctions]<0.5, p<0.01; (4) genes involved in Oct4 pathways, and (5) gene involved in stem cell pluripotency pathways using Ingunity systems software ([www.ingenuity.com](www.ingenuity.com)); (6) Expression changes of 114 epigenetic regulators between ICM and ESCs.

**Table S5** MicroRNA expression profile of ICM from E3.5 blastocysts and ESCs.

**Table S6** MicroRNA enrichment calculations using target prediction algorithms of PicTar, Miranda, or TargetScan.

**Table S7** The lists of 1,378 potential pluripotency genes (FC[Day5-Oct4$^+$/Day5-Oct4$^-$]>4, p<0.01, and r >0.6). Here r represents the averaged correlation coefficient with the expression of *Oct4*, *Sox2*, and *Nanog* in ICM outgrowth.