

Supporting Information

Wontakal et al. 10.1073/pnas.1121019109

SI Materials and Methods

ChIP-Seq Data Analysis. All reads (36 bp) from the GATA-1 ChIP-Seq were mapped to the mouse genome (version mm9 or build 37) by the ELAND aligner within the Illumina Analysis Pipeline. Aligned reads with a single best matching location and up to two mismatches were retained for peak identification. As genome-wide binding profiles from replicate chromatin immunoprecipitation and high-throughput sequencing (ChIP-Seq) runs were highly correlated (data not shown), we combined aligned reads from replicates and obtained a total of 9,486,136 and 18,717,198 uniquely mapped reads from embryonic stem cell-derived erythroid progenitor (ES-EP) proliferating and differentiating cells, respectively. For input DNA controls, 7,488,885 and 8,535,843 reads were obtained for each of these two respective conditions. The program spp (1) was applied to call peaks, genomic regions of significant GATA-1 occupancy, by comparing reads from the immunoprecipitated (IP) sample to the corresponding input DNA sample. To estimate parameters for final peak calling, we initially called peaks for replicated data independently and then the combined data and calculated the overlaps of peaks from the combined reads with the union of peaks from two individual replicates. Afterward, we selected spp parameters to maximize the overlaps that also yield comparable peak numbers between proliferating and differentiating conditions. The final spp parameters were scored >8 and $\text{enrichment.lb} >1$, with the rest set to defaults, which resulted in 6,600 and 10,600 GATA-1 peaks in proliferating and differentiating ES-EP cells, respectively.

For each GATA-1 peak we extracted a 500-bp sequence around the peak center (i.e., ± 250 bp) and used it for de novo motif discovery by the MEME software (2). The known consensus GATA-1 binding motif was the top motif returned by MEME for both proliferating and differentiating conditions, and 95% and 79% of peaks were found to contain the consensus motif, respectively.

To associate peaks with target genes, we tested several criteria, assigning genes with GATA-1 peaks from -2 kb, -10 kb, or -20 kb of transcription start sites (TSS) to $+10$ kb of transcription end sites (TES). In addition, we also tested assigning peaks to the closest gene. As a group, GATA-1 targets resulting from the first three criteria did not show significant differences of GATA-1-dependent gene expression except when they were compared with the last group on the basis of the nearest distance (Fig. S3). Accordingly, we chose -20 kb of TSS to $+10$ kb of TES as the criterion for assigning a peak to a gene. The same criterion was used to assign ChIP-Seq peaks for PU.1, Klf1, and SCL to genes. Totals of 1,380 Klf1 peaks and 2,994 SCL peaks were collected from previous studies by Tallack et al. (3) and Kassouf et al. (4) performed with fetal-liver erythroid progenitors, and PU.1 peaks (16,241) were obtained from our previous ChIP-Seq analysis in ES-EP cells (5). Gene targets were separated into groups on the basis of the co-occupancy pattern of these four transcription factors, and these groups were subjected to functional analysis with Ingenuity Pathway Analysis and DAVID GO analysis (6).

Gene Expression Analysis. Acquisition of gene expression data for proliferating and differentiating ES-EP and murine erythroleukemia (MEL) cells and fetal-liver erythroid progenitors from wild-type and PU.1 $\text{URE}^{-/-}$ was described previously (5) and can be accessed from the National Center for Biotechnology Information's Gene Expression Omnibus (GEO) database using accession no. GSE21953. Microarray data were normalized by the RMA method in the GeneSpring GX software. The \log_2

transformed signal intensities were averaged for biological replicates and used for computing expression fold change. Heat maps were generated with the mean value of all time points for a given gene and assigned a color gradient for each time point by calculating the \log_2 ratio of that time point to the mean expression value. GATA-1-dependent gene expression data were downloaded from the GEO database using accession no. GSE18042 (7). After data normalization, the fold change in expression of a gene was calculated by comparing its expression value at 0 h with the average value of all other time points. Klf1- and SCL-dependent expression data were obtained from Hodge et al. (8) and the GEO database using accession no. GSE21877 (4), respectively. The signal intensities were \log_2 transformed and quartile normalized. Fold changes in gene expression were determined with the limma algorithm (9) in the Bioconductor package.

e4c Interaction Data. A total of 551 and 273 e4C genomic clusters interacting with Hba and Hbb, respectively, and a total of 6,396 highly transcribed genes (by RNAPII-S5P occupancy) were obtained from Schoenfelder et al. (10). Active Hba- and Hbb-interacting genes in erythroid cells were defined as transcribed genes within the e4C clusters, as described previously (10). Overrepresentation of active e4C genes in different groups was calculated using the hypergeometric probability distribution.

Statistical Analysis. Wilcoxon signed rank tests were applied to compare differences in gene expression changes between any two groups of genes with designated patterns of transcriptional factor occupancy. A binomial test was applied to calculate the enrichment of genes occupied by three factors, respectively, in all genes from the mouse genome, erythroid-specific genes, and genes highly expressed in the erythroid lineage. All statistical analyses were carried out in the R language.

Mathematical Modeling. A system of four coupled nonlinear ordinary differential equations is used to model the GATA-1-PU.1 regulatory network:

$$\frac{dG}{dt} = \underbrace{G_s}_1 + \underbrace{k_{ar} \frac{G^n}{K_{ar}^n + G^n} \frac{K_{ir}^n}{K_{ir}^n + P^n}}_2 - \underbrace{k_d G}_3, \quad [S1]$$

$$\frac{dP}{dt} = P_s + k_{ar} \frac{P^n}{K_{ar}^n + P^n} \frac{K_{ir}^n}{K_{ir}^n + G^n} - k_d P, \quad [S2]$$

$$\frac{dG_T}{dt} = \underbrace{k_{at} \frac{G}{K_{at} + G}}_1 + \underbrace{k_{at} \frac{K_{it}}{K_{it} + P}}_2 - \underbrace{k_d G_T}_3, \quad [S3]$$

$$\frac{dP_T}{dt} = k_{at} \frac{P}{K_{at} + P} + k_{at} \frac{K_{it}}{K_{it} + G} - k_d P_T. \quad [S4]$$

The state variables of the preceding system of equations can be interpreted according to the following definitions: G , GATA-1 concentration; P , PU.1 concentration; G_T , GATA-1 target concentration; and P_T , PU.1 target concentration. The parameter values are labeled with subindexes a , activation; i , inhibition; d , degradation; r , regulator; t , target; s , stimulus. G_s and P_s are the

GATA-1 and PU.1 stimulation rates; n is the Hill coefficient; k_{ar} and k_{at} are maximal activation rates; K_{ar} , K_{at} , K_{ir} , and K_{it} are half-maximal concentrations for activation or inhibition as indicated by subscripts; and k_d is the first-order degradation rate assumed to be equal for all species included in the model.

To demonstrate the independence of the model from a particular choice of units for time and concentration, it is possible to introduce a change of variables that enables the expression of all parameter values as dimensionless quantities according to the Buckingham π -theorem (11). For example, time can be scaled by the degradation rate (which has units time^{-1}) rather than specifying an arbitrary scale. This process demonstrates the capacity to eliminate redundant dimensions from the parameter space of the model. We therefore derive a dimensionless form for the model of Eqs. S1–S4 in Eqs. S5–S8:

$$\frac{d\gamma}{d\tau} = \gamma_s + \kappa_r \frac{\gamma^n}{1 + \gamma^n} \frac{\lambda_r^n}{\lambda_r^n + \pi^n} - \gamma, \quad [\text{S5}]$$

$$\frac{d\pi}{d\tau} = \pi_s + \kappa_r \frac{\pi^n}{1 + \pi^n} \frac{\lambda_r^n}{\lambda_r^n + \gamma^n} - \pi, \quad [\text{S6}]$$

$$\frac{d\gamma_T}{d\tau} = \kappa_t \left(\frac{\gamma}{\alpha + \gamma} + \frac{\lambda_t}{\lambda_t + \pi} \right) - \gamma_T, \quad [\text{S7}]$$

$$\frac{d\pi_T}{d\tau} = \kappa_t \left(\frac{\pi}{\alpha + \pi} + \frac{\lambda_t}{\lambda_t + \gamma} \right) - \pi_T. \quad [\text{S8}]$$

In the derivation of the dimensionless Eqs. S5–S8 the state variables and time from Eqs. S1–S4 have been scaled according to the following relationships: $\gamma = \frac{G}{K_{ar}}$, $\pi = \frac{P}{K_{ar}}$, $\gamma_T = \frac{G_T}{K_{ar}}$, $\pi_T = \frac{P_T}{K_{ar}}$, and $\tau = k_d t$. The relationship between the parameter values of the dimensional and dimensionless forms of the model along with an associated set of base parameter values is presented in Table S1. Selecting values for the dimensionless parameters induces the definition of an equivalence class of dimensional models that all exhibit similar qualitative behavior where the relationships among the dimensioned parameter values result in the corresponding values for the dimensionless parameters. Our numerical simulations of Eqs. S1–S4 are based upon a member of the equivalence class defined by the parameter values in Table S1. Note that the parameter values K_{ir} and K_{it} (equivalently λ_r and λ_t) are modulated in numerical simulations to compare alternative network topologies (Table S2).

Here we define the assumptions of the model with reference to variables and parameter values in Eqs. S1–S4. Identical assumptions apply to Eqs. S5–S8. To model the GATA-1/PU.1 regulatory network we assume that the network architecture for GATA-1 and its targets is symmetric to that for PU.1 and its targets. Somewhat more formally, there is a permutation symmetry among the state variables representing GATA-1 and PU.1 as well as GATA-1 targets and PU.1 targets as demonstrated by the invariance of the model under the set of transformations $\sigma : \{G \leftrightarrow P, G_T \leftrightarrow P_T, G_s \leftrightarrow P_s\}$. We introduce asymmetry only in the upstream stimuli (e.g., erythropoietin and GM-CSF represented by the relative magnitudes and duration of G_s and P_s). In Eqs. S1 and S2, G_s and P_s , respectively represent GATA-1 and PU.1 upstream stimulation rates. The second terms of Eqs. S1 and S2 consist of three components. The first is the maximal activation rate described by the parameter k_{ar} . The second parts are Hill functions describing the autoregulation of GATA-1 and PU.1 with half-maximal activation constants K_{ar} (12). The corresponding Hill coefficients, n , in the base parameter set are >1 to represent the existence of multiple binding sites for GATA-1 and PU.1 in the upstream regulatory regions of the GATA-1 and

PU.1 genes (13–15). The third parts are Hill functions representing the mutual inhibition of PU.1 and GATA-1 on each other's gene expression with half-maximal inhibition constants K_{ir} . Qualitatively identical results are obtained even if the Hill coefficients in the autoregulatory and cross-inhibition terms are independent for all combinations of Hill coefficients in the range we tested: $n = 2, \dots, 6$. The autoregulatory and cross-inhibition components of the second terms are multiplied by one another to represent their competition to control the synthesis rates of GATA-1 and PU.1. The final terms of Eqs. S1 and S2 represent first-order degradation processes with rates k_d for both GATA-1 and PU.1. Eqs. S3 and S4 represent the dynamics of GATA-1 and PU.1 targets, respectively. The first term of Eq. S3 represents GATA-1-mediated activation of its targets with half-maximal activation constant K_{at} , the second term is PU.1 inhibition of the expression of GATA-1 targets with half-maximal inhibition constant K_{it} , and the third term is the first-order degradation with rate k_d of the GATA-1 targets. We have assumed that GATA-1 and PU.1 serve as independent inputs to their respective target genes. The terms of Eq. S4 are analogous to those of Eq. S3 but describe the regulation of PU.1 targets. The overall form of Eqs. S1–S4 is similar to that proposed by Laslo et al. to model a different aspect of the hematopoietic gene regulatory network (16). This system of equations is symmetric for GATA-1 and PU.1 and therefore the differential stimulus applied to favor GATA-1 and the erythroid cell fate over PU.1 and the myeloid cell fate in the test case shown in Fig. 4 would produce precisely the opposite result (i.e., PU.1-mediated myeloid lineage differentiation as opposed to GATA-1-mediated erythroid lineage differentiation were the stimuli magnitudes permuted).

Characterization of the Mathematical Model. To understand the function of the GATA-1–PU.1 network topology identified experimentally (Fig. S7D) we investigated the effects of continuous modulation of the network topology upon the cell fate determination specified by the ratio between the steady-state expression levels of GATA-1 and PU.1 target genes. The model was evaluated via numerical simulations in which GATA-1 receives an initial transient stimulus G_s that is 10-fold higher than that applied to PU.1 (Figs. S8–S11 and Table S1). To produce the network topology shown in Fig. S7A given the system of Eqs. S1–S4 we require the parameters K_{ir} and K_{it} to take values that are high relative to the maximal protein concentrations. To modulate the network topology from that of Fig. S7A (bottom center corner of Fig. 4A) to that of Fig. S7B (right corner of Fig. 4A) we varied the parameter K_{ir} along this axis (the x axis) from high [$\text{Max}(K) = 10^3$] to low [$\text{Min}(K) = 10^{-1}$]. However, for clarity of presentation, on the x axis of Fig. 4 we transformed the K_{ir} values using the following function $f(K) = \text{Max}(K) - \text{Min}(K) - K_j$, where K represents the vector of K_{ir} values $K = \{K_j\}$, thus representing the antagonistic interaction *strength*. We transformed the K_{it} values in the same way to represent antagonistic interaction strengths along the y axes of Fig. 4. When K_{ir} is high, the terms $\frac{K_{ir}^n}{K_{ir}^n + P^n}$ and $\frac{K_{it}^n}{K_{it}^n + G^n}$ from Eqs. S1 and S2, respectively, are ≈ 1 and therefore neither GATA-1 nor PU.1 inhibits the expression of the other. As K_{ir} decreases, the concentrations of GATA-1 and PU.1 play increasingly significant roles as inhibitors of the expression of the other and, in the case when either GATA-1 or PU.1 reaches extremely high levels, these terms approach zero. The network topology is similarly modulated from that of Fig. S7A to that of Fig. S7C by decreasing the value of the parameter K_{it} . The network topology represented in Fig. S7D is produced when both K_{ir} and K_{it} take on low values relative to the maximal protein concentrations.

To characterize the dynamics of GATA-1–PU.1 regulation near each of the four corners of the $K_{ir} - K_{it}$ parameter space represented in the xy plane of Fig. 4A we simulated Eqs. S1–S4 for four

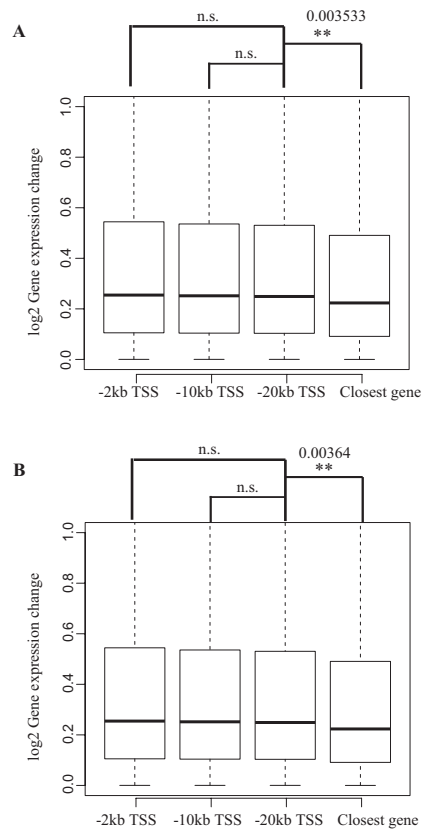


Fig. S3. Analysis of four different criteria for assigning GATA-1 ChIP-Seq peaks to genes. GATA-1 ChIP-Seq peaks in proliferating ES-EP (A) and differentiating ES-EP (B) were assigned to genes on the basis of four different criteria: (i–iii) assignment if the peak lies within the region spanning from –2 kb (i), –10 kb (ii), and –20 kb (iii) of a TSS through +10 kb of the TES and (iv) assignment of the peak to the nearest gene. Using these four assignment criteria and published GATA-1–dependent gene expression changes in erythroid cells (1, 2), boxplots were constructed to show the log₂ fold change in expression of the four different sets of genes.

- Cheng Y, et al. (2009) Erythroid GATA1 function revealed by genome-wide analysis of transcription factor occupancy, histone modifications, and mRNA expression. *Genome Res* 19: 2172–2184.
- Yu M, et al. (2009) Insights into GATA-1-mediated gene activation versus repression via genome-wide chromatin occupancy analysis. *Mol Cell* 36:682–695.

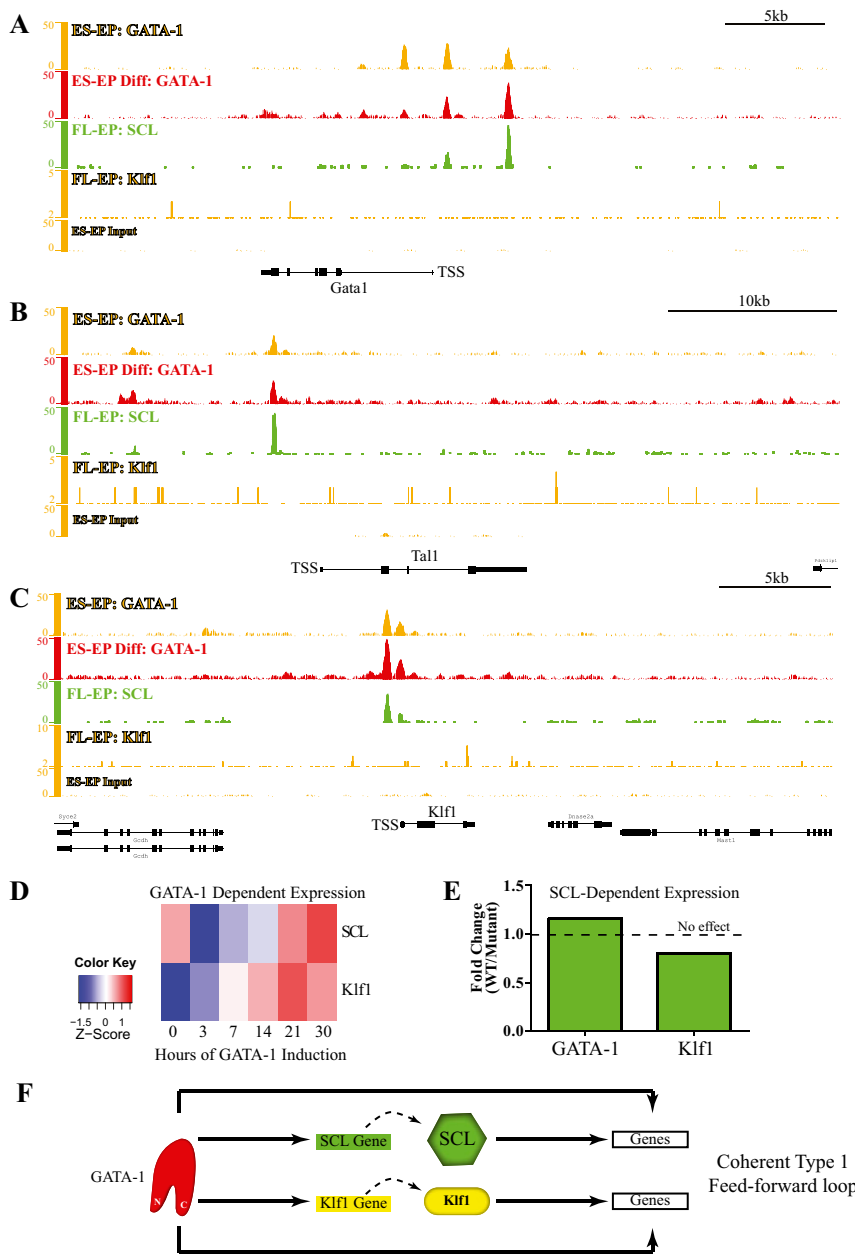


Fig. S4. GATA-1 acts upstream of SCL and Klf1. (A–C) Occupancy maps from ChIP-Seq data for GATA-1 in proliferating and differentiating ES-EP and for SCL and Klf1 in FL-EP are shown in the vicinity of the genes encoding GATA-1 (A), SCL (Tal1) (B), and Klf1 (C). (D) Heatmap of GATA-1–dependent changes in gene expression (1, 2) of SCL and Klf1. (E) Gene expression differences for GATA-1 and Klf1 between FL-EP from wild-type and SCL DNA-binding–defective mutant mice (3). (F) A model for a coherent type 1 feed-forward loop formed by GATA-1 and SCL and by GATA-1 and Klf1 based on the ChIP-Seq and gene expression data displayed in A–E.

- Cheng Y, et al. (2009) Erythroid GATA1 function revealed by genome-wide analysis of transcription factor occupancy, histone modifications, and mRNA expression. *Genome Res* 19: 2172–2184.
- Yu M, et al. (2009) Insights into GATA-1-mediated gene activation versus repression via genome-wide chromatin occupancy analysis. *Mol Cell* 36:682–695.
- Kassouf MT, et al. (2010) Genome-wide identification of TAL1's functional targets: Insights into its mechanisms of action in primary erythroid cells. *Genome Res* 20:1064–1083.

Table S2. Correspondence between network topologies and parameter values

Fig. 4A corner	Fig. S7	Figure containing dynamics	K_{ir}	K_{it}	Direct cross-inhibition strength	Downstream target inhibition strength
Center bottom	A	Fig. S8	100	100	Low	Low
Right	B	Fig. S9	1	100	High	Low
Left	C	Fig. S10	100	1	Low	High
Center top	D	Fig. S11	1	1	High	High

Table S3. List of qChIP primers

Gene	Forward primer	Reverse primer	Reference
<i>Myogenin</i>	GAA TCA CAT GTA ATC CAC TGG A	ACG CCA ACT GCT GGG TGC CA	(1)
<i>β-HS2</i>	TGT GTT CAG CCT TGT GAG CCA GC	TGG ACT TCC TCC TAG AGA CCC AG	(2)
<i>HDGF</i>	CCA AGA AAG ATG TGG GAG GA	CTG CTG CAG AAA GCT GAT TG	This study
<i>Zdhhc19</i>	TTT GAG GGT GAG GGT CAA AG	CCA TTT CTG CCA GGA GGT TA	This study
<i>Slc16a10</i>	CTG CAG AGG CCA GAT AAG GA	AAG CTA GGG GAC AAG GGA TG	This study
<i>Gypc</i>	CAC GCC TAT CAG CAT ATG GA	GAG ACA GCT ACC ACG GGT GT	This study
<i>Chr 13 140613517</i>	CAG GCT GGG AGA GAA TTT TG	GTA CGC ACT TTG GGG TTT TG	This study
<i>H2AFY3</i>	GGT CCA GGA CAA CGG TTC TA	AGC TCA GGG TGT GAC AGA GG	This study
<i>Ifih1</i>	CCC TTA TCA ATG GCC ACA GT	AAA ACG GAA TCA ACG GTT TG	This study
<i>Dapp1</i>	GCC AAT GCA TAA GTG AGC AA	GGC TTC CGG AAC ACA AGA TA	This study
<i>Accn2</i>	AAT CGG AAA GAT CCC AGC TT	AAT GCA GCC CTC CAT ATC AC	This study
<i>Gfi-1b</i>	GCC CCT GAT AAC ACT TGG AA	GCA ACT GGA GGG AAA TCT GA	This study
<i>Fam125B</i>	TAT GTC TGG TGG CAC ATG CT	GTG ACA GCC AAA GGA GGA AA	This study
<i>Zfpm1</i>	AGC GAT GGG GTT GAT AAG TG	CGG TGA TAA GCA GAG CCT GT	This study
<i>Lyl1</i>	GGG GTC AGC ATT GCT TCT TA	CCT GGC TTC CTC CCT CTT AC	This study
<i>Chr 16 93147915</i>	TAC CCT GGT CTC ACC TCA GC	AGG CAG TGA AGG GGA AAG AT	This study
<i>Pvt1</i>	GAT GTC CCC AGA TAG CCA GA	AGA CTC CAG AAG TGG GCT GA	This study
<i>St3gal</i>	TTG CGA ACA TGC AAA GAT TC	TTT GAG AAG AGT GGG GCA GT	This study
<i>Rap1a</i>	GGT CGT TTG TCT TTC CTC CA	CAG TCT GTC CCC TCC CTA CA	This study
<i>Trem14</i>	GGC CTG CTG AAT TCT CTC TG	GGA GAA GGA ACC TCC TGA CC	This study

- Lu J, McKinsey TA, Zhang CL, Olson EN (2000) Regulation of skeletal myogenesis by association of the MEF2 transcription factor with class II histone deacetylases. *Mol Cell* 6:233–244.
- Stopka T, Amanatullah DF, Papetti M, Skoultchi AI (2005) PU.1 inhibits the erythroid program by binding to GATA-1 on DNA and creating a repressive chromatin structure. *EMBO J* 24: 3712–3723.