# Regulation of genome-wide transcription by essential factors that control promoter-proximal RNA polymerase II pausing in human cells.

Student: *Mariia Vlasenok*

Research Advisors: *Konstantin Severinov*

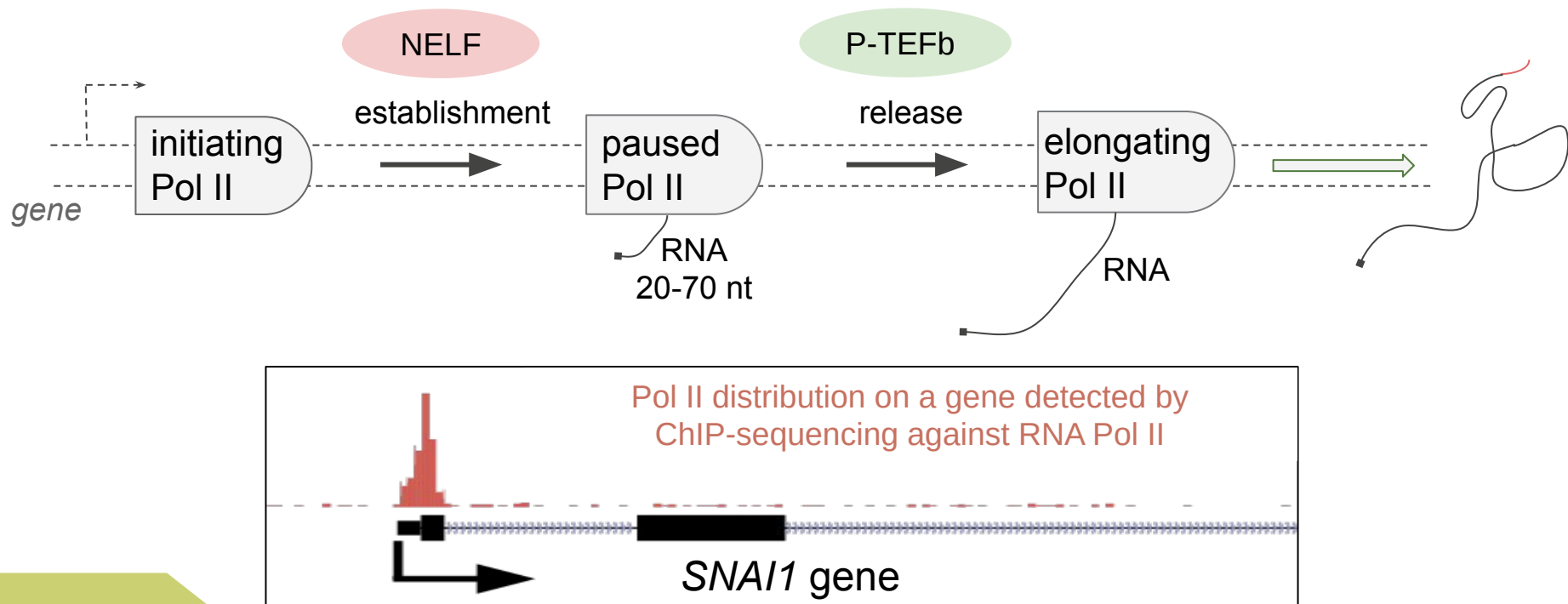*Sergei Nechaev*

*Yen Lee Loh*

June, 2019

# Background

- In eukaryotes, messenger RNA genes are transcribed by the RNA Polymerase II (Pol II).

- The main stages of Pol II transcription are initiation, elongation, and termination.

- Pol II transcription elongation is regulated as tightly as initiation.

# Background

- Promoter-proximal Pol II pausing - halt of Pol II less than 100 nucleotides downstream of the start site.

- NELF and P-TEFb control, respectively, pausing establishment and release on every gene.



Pol II distribution on a gene detected by ChIP-sequencing against RNA Pol II

*SNAI1* gene

# Background

- Pol II pausing is associated with <u>active,</u> not repressed genes.

- Pausing (and NELF + P-TEFb) is <u>essential</u> in higher eukaryotes.

- Pausing is implicated in organism development, cell differentiation and responses to stimuli, but the mechanisms are unknown.

**How can <u>essential</u> factors <u>regulate</u> transcription?**

# Background

**How can <u>essential</u> factors <u>regulate</u> transcription?**

**Overall idea:**

Pausing regulates the distribution of Pol II on gene promoters across the genome.

**Working hypothesis:**

Limiting the levels of NELF and P-TEFb favors Pol II transcription at stronger promoters at the expense of less active ones.

**Significance:**

Limiting the levels of essential factors may be a common mechanism organizing transcription into stable genome-wide patterns.

# Aims of the project

## Aim 1

to analyze transcriptional responses of human cells to heat shock using global run-on sequencing data

*bioinformatics*

## Aim 2

to model transcriptional consequences of NELF depletion in a cell as a closed system
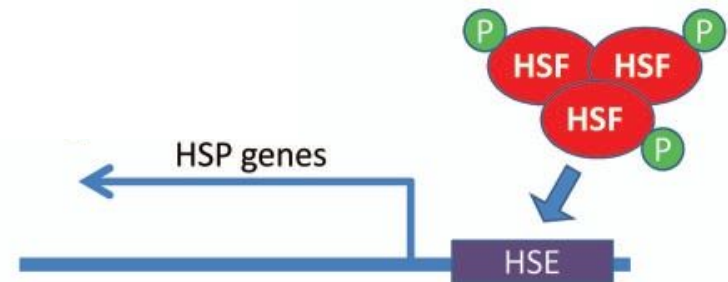
*mathematical modeling*

# Aim 1

## to analyze transcriptional responses of human cells to heat shock using global run-on sequencing data

to compare Heat Shock (HS) response
in completely different human cells

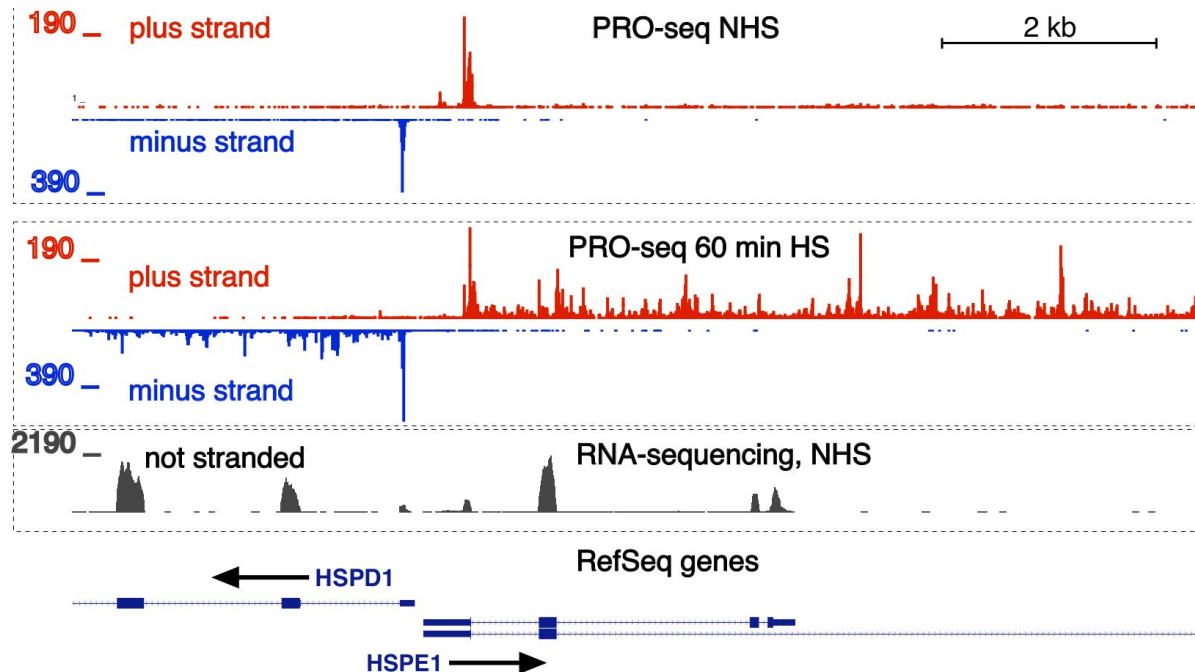- K562 - leukemia cell line
- MCF-7 - breast cancer cell line

Heat Shock response includes
conserved activation
by Heat Shock Response factor (HSF)

# Aim 1 methods

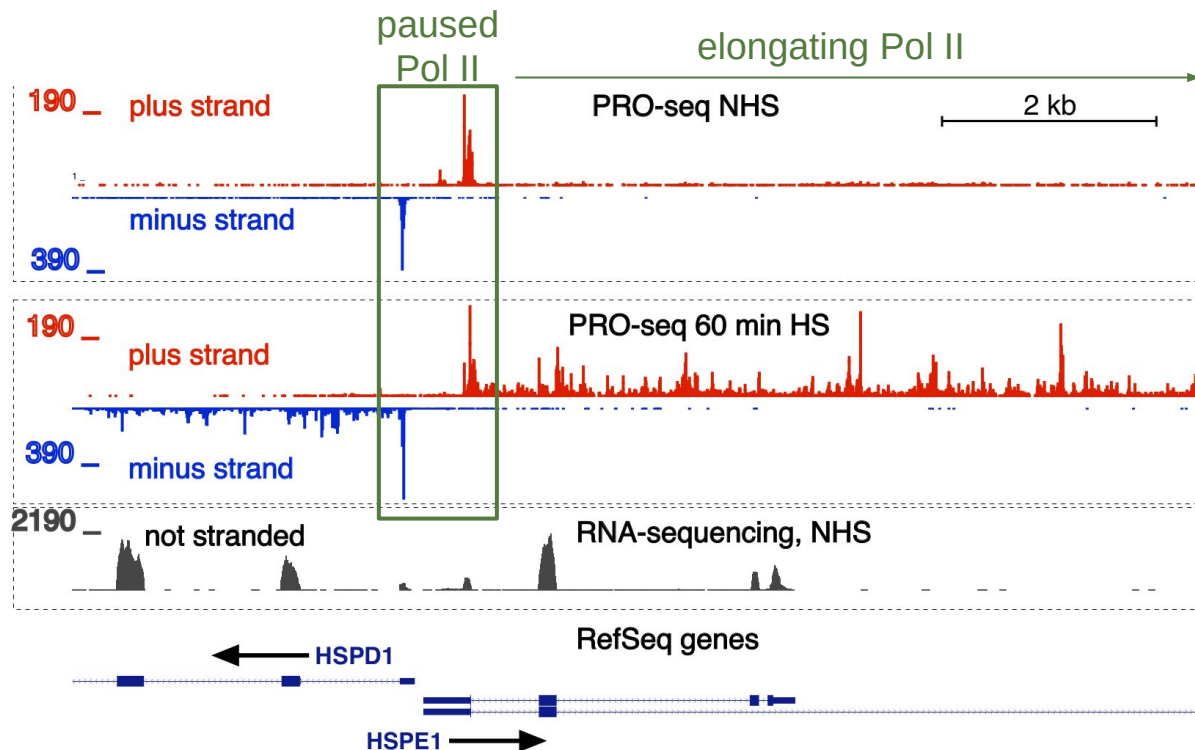PRO-seq - Precision nuclear Run-On sequencing

- uses biotin labeled nucleotides to detect <u>nascent transcription</u>
- enables genome-wide mapping of transcriptionally engaged

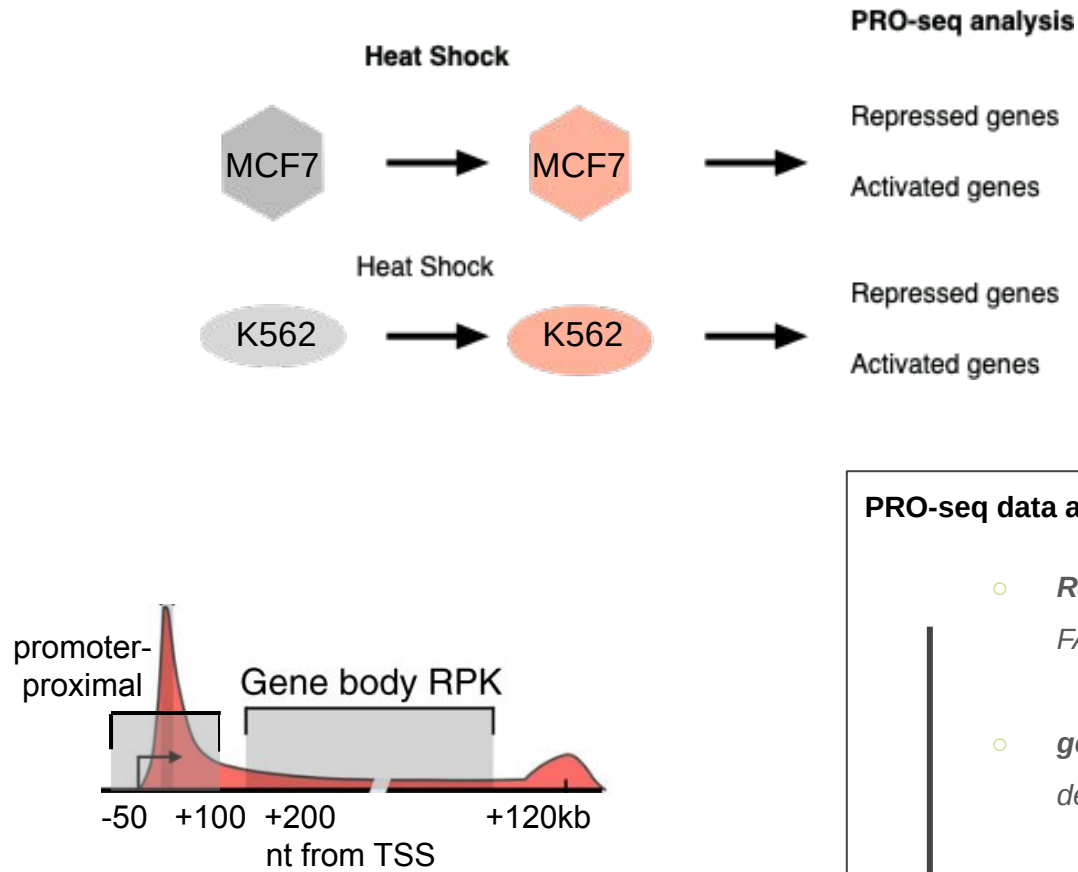  RNA Polymerase with single nucleotide-resolution

# Aim 1 methods

PRO-seq - Precision nuclear Run-On sequencing

- uses biotin labeled nucleotides to detect <u>nascent transcription</u>
- enables genome-wide mapping of transcriptionally engaged RNA Polymerase with single nucleotide-resolution
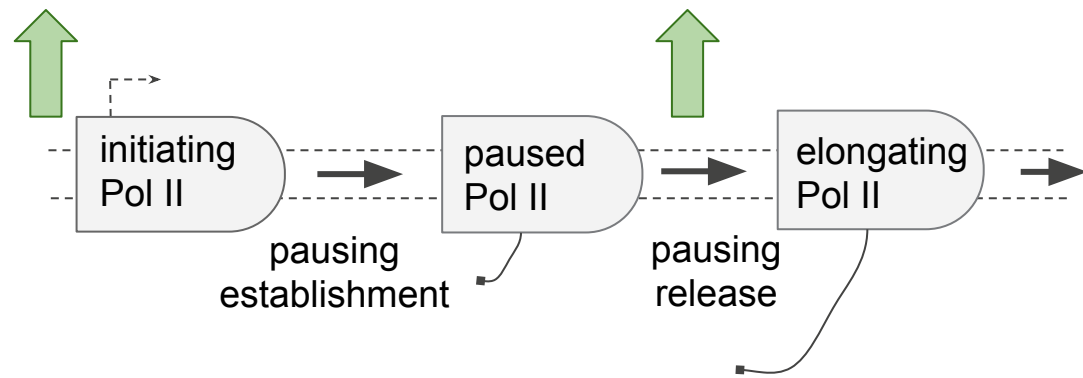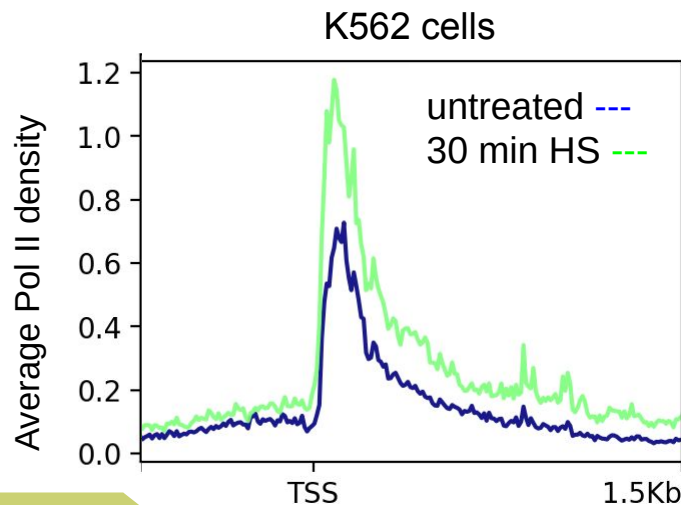
# Aim 1 methods

**Heat Shock**

MCF7 → MCF7 →

**PRO-seq analysis**

Repressed genes

Activated genes

**Heat Shock**

K562 → K562 →

Repressed genes

Activated genes

promoter-proximal
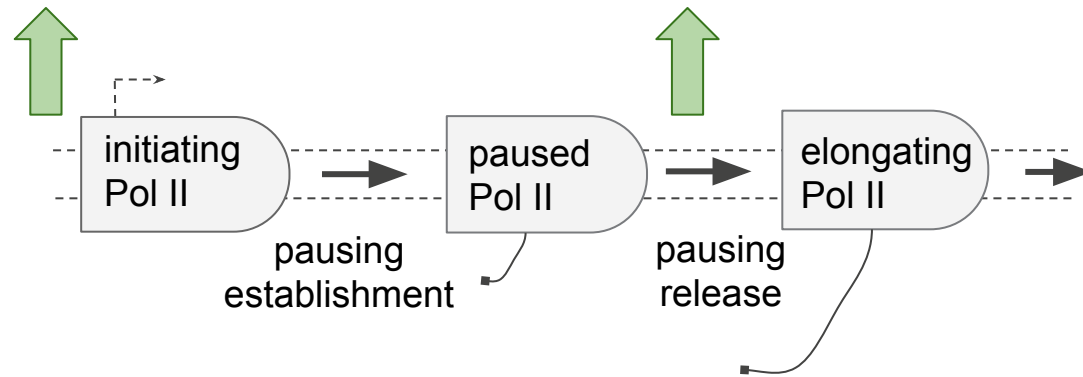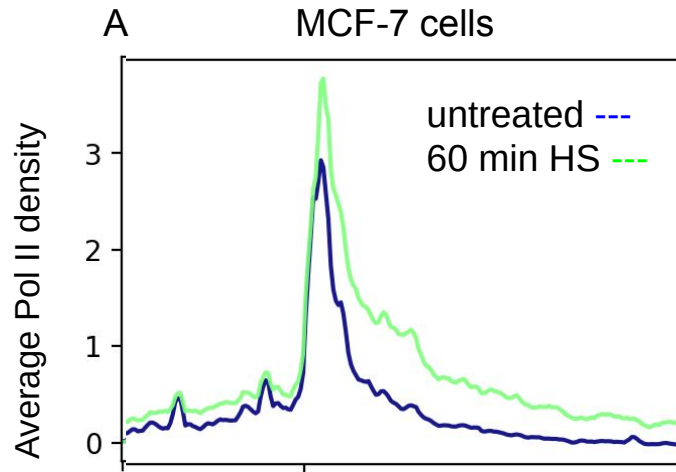
Gene body RPK

-50   +100   +200   +120kb

nt from TSS

**PRO-seq data analysis:**

○ *Raw data: FastQC, FASTX-toolkit, hisat2, samtools*

○ *genome arithmetic - bedtools, deeptools*

○ *Normalization*

○ *Differential expression: DESeq2*

# Heat Shock response: similar mechanisms of activation

in two cell lines

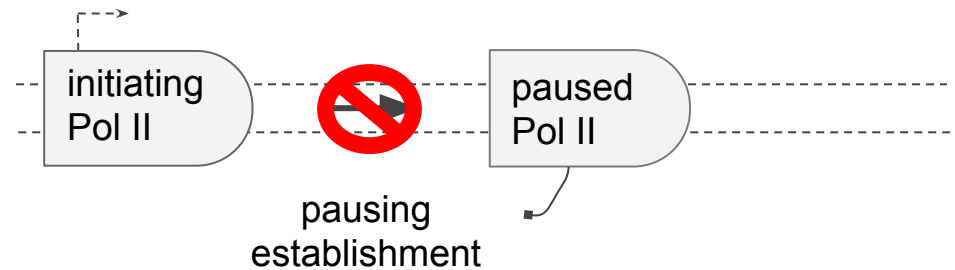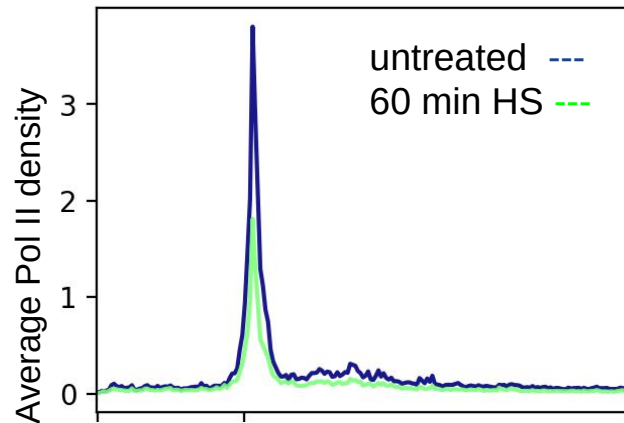Differentially expressed genes (p adj < 0.05) <u>activated</u> during HS.

# HS response: different mechanisms of repression

in two cell lines

Genes transcriptionally (p adj < 0.05) <u>repressed</u> during HS.

MCF-7 cells



K562 cells



This mechanism of regulation without initiation repression was observed in both previously studied with pro-seq cell types (MEFs K562)

# Aim 1 Conclusions

- We compared transcriptional response of distinct two human cell lines to the same stimulus.

- We find major differences in mechanisms of gene repression. We show that repression can take place either at the level of Pol II recruitment to promoters or Pol II release from pausing.

- **SIGNIFICANCE.** The work provides evidence for transcription regulation by distribution of essential factors across active promoters.

# Aim 2

## to model transcriptional consequences of NELF depletion in a cell as a closed system



factor essential for transcription (NELF)

gene

transcription machinery

NELF ↓

cell nucleus

We suggest that the effects of NELF depletion in the cell may be explained by competition of genes for the available essential factor. The number of genes that can be active is limited by the amount of the essential factor.

# The chain of reactions used for simulation



*for each gene g:*

$$\begin{cases} P_g + e \leftrightarrow PIC_g \\ PIC_g + n \rightarrow PEC_g \\ PEC_g \rightarrow P_g + e + n \end{cases}$$

*P* - gene promoter
*e* - Pol II
*n* - NELF
*PIC* - pre-initiation complex
*PEC* - paused elongation complex
*g* - gene ( out of 1000 )

$k^i$ - initiation frequency  i.e.
        ~ promoter strength
$k^p$ - pausing frequency i.e.
        ~ NELF preferences

# Reaction Rate Equations Approach

*RRE works for systems with large numbers of molecules of each species, which is clearly not true for promoters of each gene ($p_g$). So, we model many cells (nuclei) simultaneously as if they shared all genes and factors.*

*assumption*

# Reaction Rate Equations approach

*steady state simulation*

*RRE allows to calculate concentrations in equilibrium:*

$$\frac{d}{dt}[p_g] = \frac{d}{dt}[PIC_g] = \frac{d}{dt}[PEC_g] = \frac{d}{dt}[e] = \frac{d}{dt}[n] = 0$$

Results in:

$$\begin{cases} for \ \forall \ g \ < G \\[6pt] [e]\dfrac{k_g^i}{k^t + [n]k_g^p} = \dfrac{[PIC_g]}{[P_g]} \\[10pt] [n]\dfrac{k_g^p}{k_g^r} = \dfrac{[PEC_g]}{[PIC_g]} \\[10pt] [PEC_g] + [PIC_g] + [P_g] = [P_g]_0 \\[6pt] \text{and for the whole system:} \\[6pt] [n] + \sum_g [PEC_g] = [n]_0 \\[6pt] [e] + \sum_g ([PEC_g] + [PIC_g]) = [e]_0 \end{cases}$$

We also can use the initial conditions:

*for each gene g*

# Results of simulation:
# nonlinear transcriptional response to depletion of NELF

*The simulation prediction: most active genes are activated, but the rest - repressed.*



Rank-size distribution of genes based on simulated data

# Aim 2 Conclusions

- Limiting the levels of a pausing factor NELF may be sufficient to enforce non-linearity in activity of promoters.

- **SIGNIFICANCE.** Limitation of factors at distinct steps of Pol II pausing may be a universal mechanism that stabilizes transcriptomes in different metazoan cell types.

# My work thus far　　　Future directions

# Acknowledgements

**Research advisors:**

- *Konstantin Severinov, Skoltech*
- *Sergei Nechaev,*
  *University of North Dakota School of Medicine*
- *Yen Lee Loh, University of North Dakota*

**Nechaev lab members:**

- *Nii Koney-Kwaku Koney*
- *Sayantani Ghosh Dastidar*
- *Bo Lauckner*

# Thank you!

Student: *Mariia Vlasenok*

# Heat-induced repression is much more extensive in K562 than in MCF-7 cells

Differentially expressed genes

**Activated genes**

**Repressed genes**



339 | 479 | 1285

337 | 430 | 6276

MCF-7 cells          K562 cells

## Reaction Rate Equations Approach

- **strengths**
  - "textbook" model for chemical kinetics
  - the least complex one that can account for thousands of genes
  - convenience when working with systems in equilibrium state
- **weakness**

  RRE works for systems with large numbers of molecules of each species, which is clearly not true for promoters of each gene ($p_g$).

## Chemical Master Equations Approach

- **strengths**
  - cellular processes involve extremely small population sizes, where it is unrealistic to think in terms of concentration
  - vital when system exhibits bistability
- **weakness**
  - computational complexity

# Imply constants $k^p$ and $k^i$ from experimental data

experimental data:

$$[PIC_g] + [PEC_g] = promoter\ occupancy$$

$$r_{release} = [PEC_g]k_g^r = [n]k_g^p[PIC]_g = expression$$

initial system:

$$(k^t + [n]k_g^p)[PIC_g] = e[P_g]k_g^i$$

$$[n]k_g^p[PIC]_g = [PEC_g]k_g^r$$

initial conditions:

$$[PEC_g] + [PIC_g] + [P_g] = [P_g]_0$$

$$[n] + \sum_g [PEC_g] = [n]_0$$

$$[e] + \sum_g ([PEC_g] + [PIC_g]) = [e]_0$$

$k^p$, $k^i$ and $k^t$ don't change after NELF depletion; $n$, $n_0$ and $e$ do.

System with *2g+2* variables

$$g+2 \begin{cases} expr_g(nk_g^p + k^t) = nk_g^p(cells - prom_g)ek_g^i \\ e + \sum prom_g = e_0 \\ n + \sum prom_g - \sum \frac{expr_g}{nk_g^p} = n_0 \end{cases}$$

*2g+4* variables and *2(g+2)* equations with NELF‑‑ data.

rank-size distribution of genes



simulated $k^i$, calculated $nk^p$

genes, ordered by $k^p$

# Appendix
*pipeline*

1. trim adaptor (3' TGGAATTCTCGGGTGCCAAGG) and trim reads to min length of 15bp
2. exclude reads that map to ribosomal RNA (*hisat2)*
3. alignment to hg19 using *hisat2* (uniquely aligned reads with no more than 2 mismatches)
4. read counts for gene bodies
   a. 23k genelist
   b. gene body is considered to start at TSS+200bp and end at TSS+120kb or gene end, whichever is closer
   c. *bedtools intersect* were used to count overlaps
5. Normalization using 3' ends of long genes
   a. used genes > 150kb
   b. normalization region: from TSS + 100kb to TTS - 0.5kb
6. DESeq2 on gene body read densities (reads per 10kb)
   a. p-value < 0.05
   b. no FC cutoff



>150 kb genes

Normalization region

polyA

30 min * 3kb/min = 90kb    +100 kb from TSS    -0.5 kb from polyA

*Vihervaara et al, 2017.*

# Heat Shock response includes conserved activation and variable repression



**MA-plots:**

**MCF-7 cells**
PRO-seq change 30HS:NHS
gene body := [TSS+200,TSS+120kb]

**K562 cells**
PRO-seq change 30HS:NHS
gene body := [TSS+200,TSS+120kb]

$p < 0,05$

activated

repressed

activated

repressed

log fold change upon HS

Mean of normalized counts in gene body, RPK

**Activated genes**

339 | 479 | 1285

MCF-7 cells

K562 cells

**Repressed genes**

337 | 430 | 6276

MCF-7 cells

K562 cells

# Repressed genes differ between cell lines

**MA-plots:**

**MCF-7 cells**
PRO-seq change 30HS:NHS
gene body := [TSS+200,TSS+120kb]

**K562 cells\***
PRO-seq change 30HS:NHS
gene body := [TSS+200,TSS+120kb]

$p < 0,05$

activated

repressed

log fold change upon HS

activated

repressed

Mean of normalized counts in gene body, RPK

**Activated genes**

483 | 448 | 894

in MCF-7 cells
in K562 cells

**Repressed genes**

552 | 1175 | 5204

in MCF-7 cells

in K562 cells

*\* K562 data re-analyzed from Vihervaara et al, 2017. Analysis results are consistent.*

## Reaction Rate Equations Approach

$$\begin{cases} P_g + e \leftrightarrow PIC_g \\ PIC_g + n \rightarrow PEC_g \\ PEC_g \rightarrow P_g + e + n \end{cases}$$

By the law of mass action

$$\begin{cases} for \; \forall \; g \; < G: \\ \frac{d}{dt}[p_g] = -k_g^i[p_g][e] + k^t[PIC_g] + k_g^r[PEC_g] \\ \frac{d}{dt}[PIC_g] = k_g^i[p_g][e] - k^t[PIC_g] - k_g^p[PIC_g][n] \\ \frac{d}{dt}[PEC_g] = k_g^p[PIC_g][n] - k_g^r[PEC_g] \\ \text{and for the whole system:} \\ \frac{d}{dt}[e] = \sum_g (-k_g^i[p_g][e] + k^t[PIC_g] + k_g^r[PEC_g]) \\ \frac{d}{dt}[n] = \sum_g (-k_g^p[n][PIC_g] + k_g^r[PEC_g]) \end{cases}$$

MA-plots:

MCF-7 cells
PRO-seq change As:Untr

K562 cells
PRO-seq change As:Untr

$p < 0,05$

log fold change after treatment

activated

repressed

activated

repressed

Mean of normalized counts in gene body, RPK

**Activated genes**
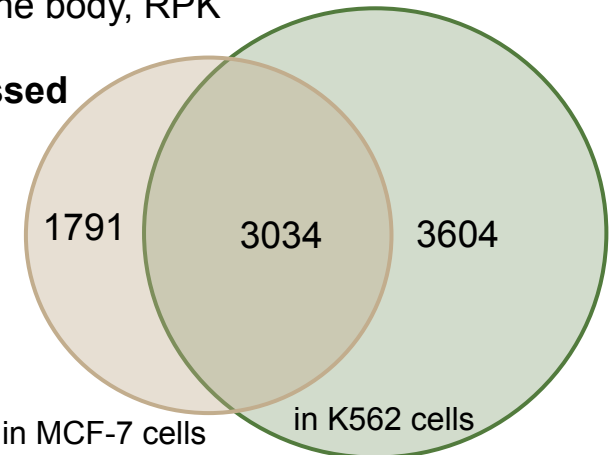
1514  1458  1180

in MCF-7 cells

in K562 cells

**Repressed genes**

1791  3034  3604

in MCF-7 cells

in K562 cells

# PRO-seq

*map the location of active RNA polymerases (PRO-seq)*

- *genome-wide*
- *strand-specific*
- *single nucleotide-resolution*

*It is used for studying short-time transcriptional responses.*

*Nuclei are isolated from cells and in vitro transcriptionally engaged RNA polymerases incorporate one biotin-NTPs into the 3′ end of nascent RNA. The biotin-labeled nascent RNA is used to prepare sequencing libraries, which are sequenced from the 3′ end to provide high-resolution positional information for the RNA polymerases.*



Cell culture

Sample preparation (nuclei isolation/cell permeabilization)

Nuclear run-on and RNA isolation

RNA fragmentation

Biotin RNA enrichment

3′ RNA adaptor ligation and second biotin RNA enrichment

Enzymatic modification of the RNA 5′ ends

5′ RNA adaptor ligation and third biotin RNA enrichment

Reverse transcription

PCR amplification and library size selection

Sequencing and data analysis

RNAP
Nascent RNA

PRO-seq

5′ De-capping

5′ Phosphorylation

From the 3′ of RNA

- ○ 5′ cap
- ● Biotin-NTP
- P- 5′ phosphate
- —— VRA3(RA5) RNA adaptor
- —— VRA5(RA3) RNA adaptor
- —— RP1 RT primer
- —— RTP RT primer
- —— RPI-n indexed primer
- ···▶ Sequencing direction