



**Skoltech**

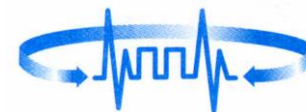
Сколковский институт науки и технологий



Учебно-Научный Центр



Биоинформатика



# Биология больших данных

*Михаил Гельфанд*

JPoint, 6 IV 2018

# Краткий курс молекулярной биологии

- Центральная догма
  - Репликация, транскрипция, трансляция
- Регуляция
  - Сигнальные пути
  - Факторы транскрипции
  - Эпигенетика
  - Пространственная структура ДНК

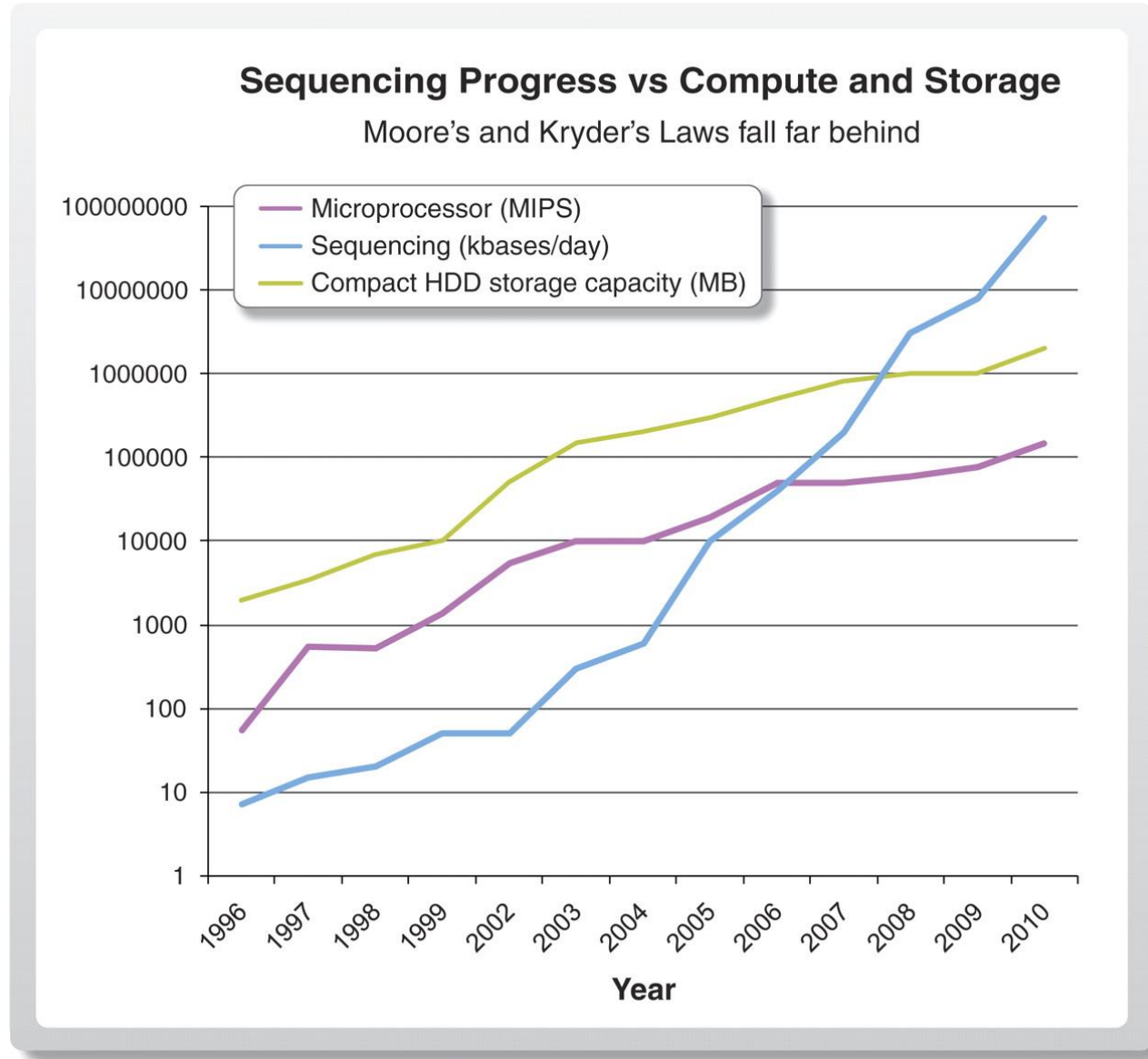
# Данные

- Уровень экспрессии
  - Концентрации мРНК
  - Концентрации белков
  - Время жизни мРНК и белков
  - Концентрации метаболитов
- Взаимодействия
  - Белок-ДНКовые
  - Белок-белковые
- Структура генома
  - Метилирование и открытость ДНК
  - Положение нуклеосом и модификация гистонов
  - Пространственная структура
- Функционально-генетические
  - Летальность и фенотип мутаций
  - Синтетические летали

# Многие методы основаны на секвенировании

- Уровень экспрессии
  - Концентрации мРНК:  
секвенирование транскриптома
- Взаимодействия
  - Белок-ДНКовые:  
ChIP-Seq
- Структура генома
  - Метилирование ДНК:  
бисульфитное секвенирование
  - Положение и модификация нуклеосом:  
ChIP-Seq
  - Пространственная структура:  
HiC

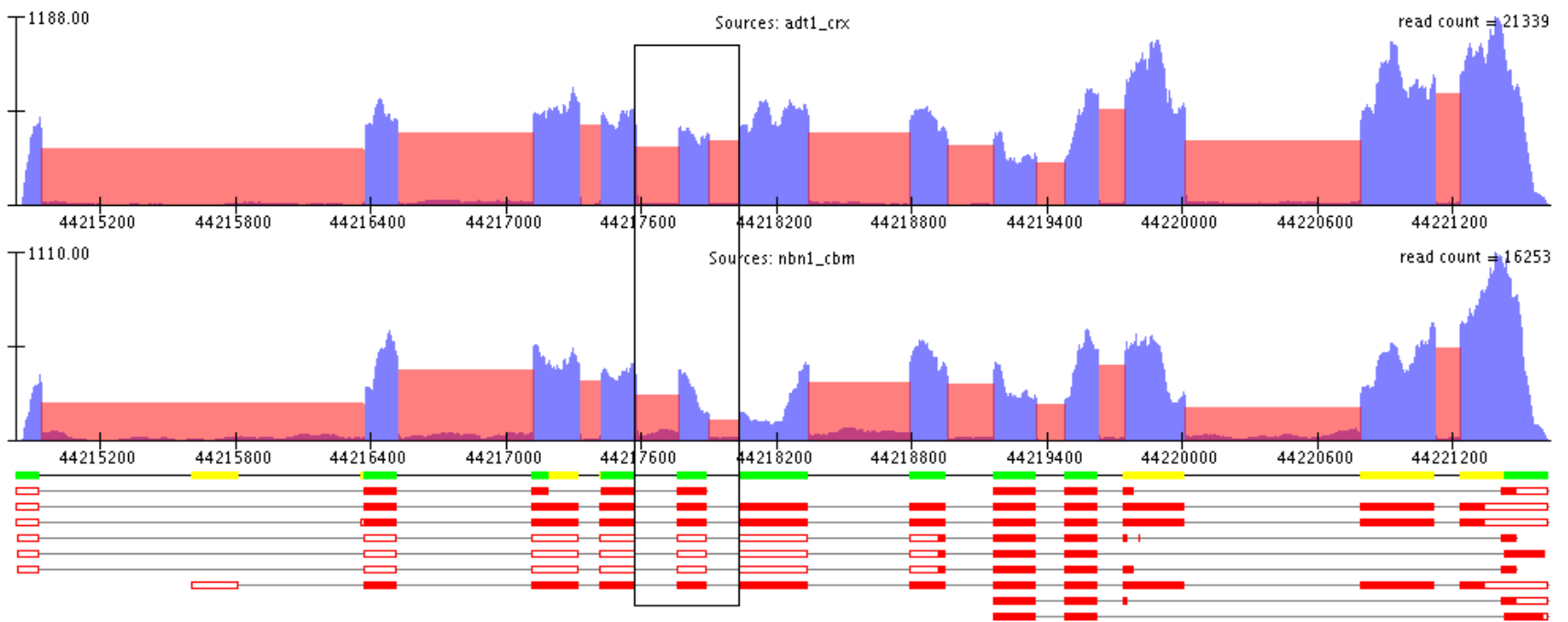
# A doubling of sequencing output every 9 months has outpaced and overtaken performance improvements within the disk storage and high-performance computation fields



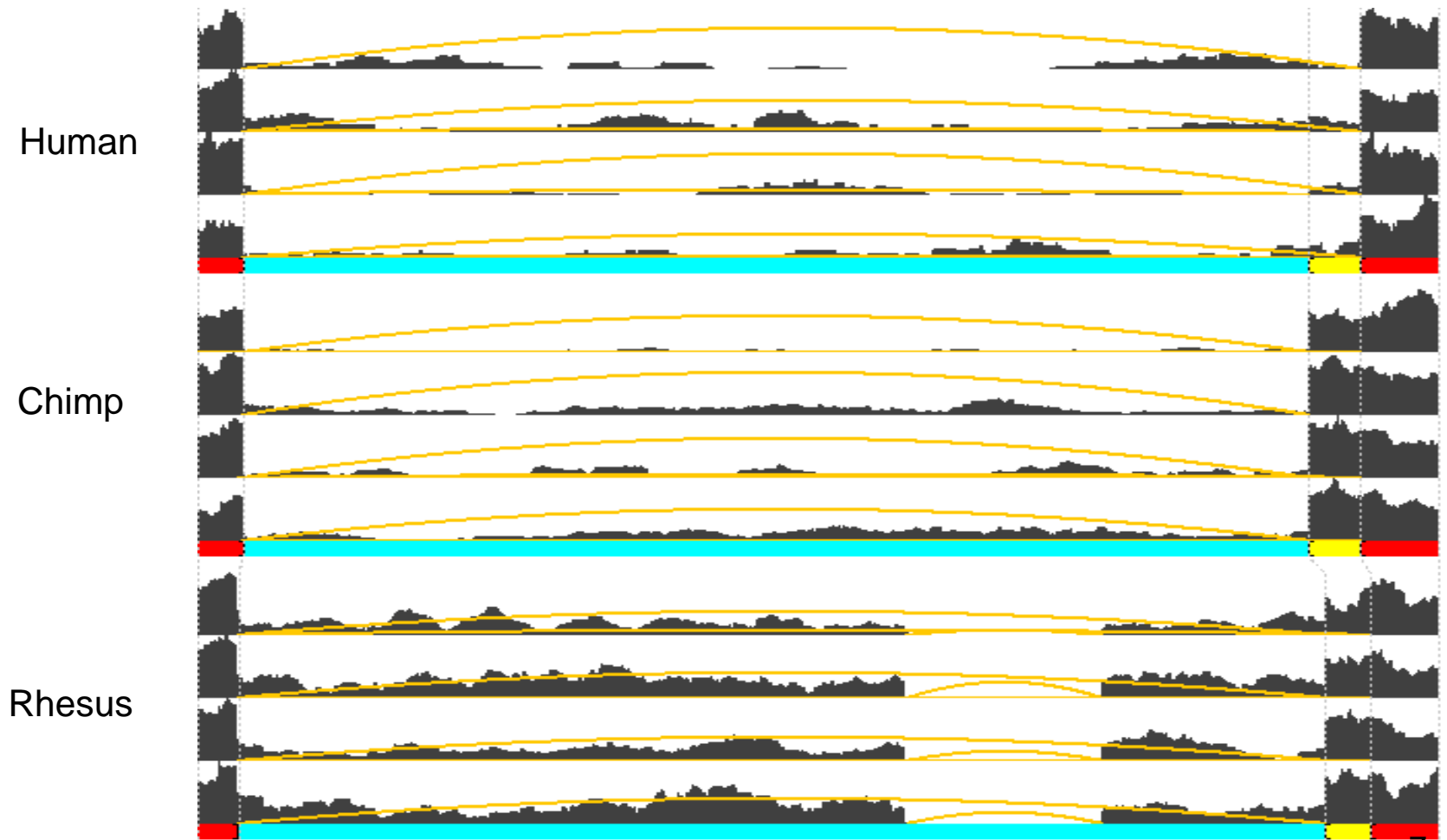
S D Kahn Science  
2011;331:728-729



# Картирование транскриптов

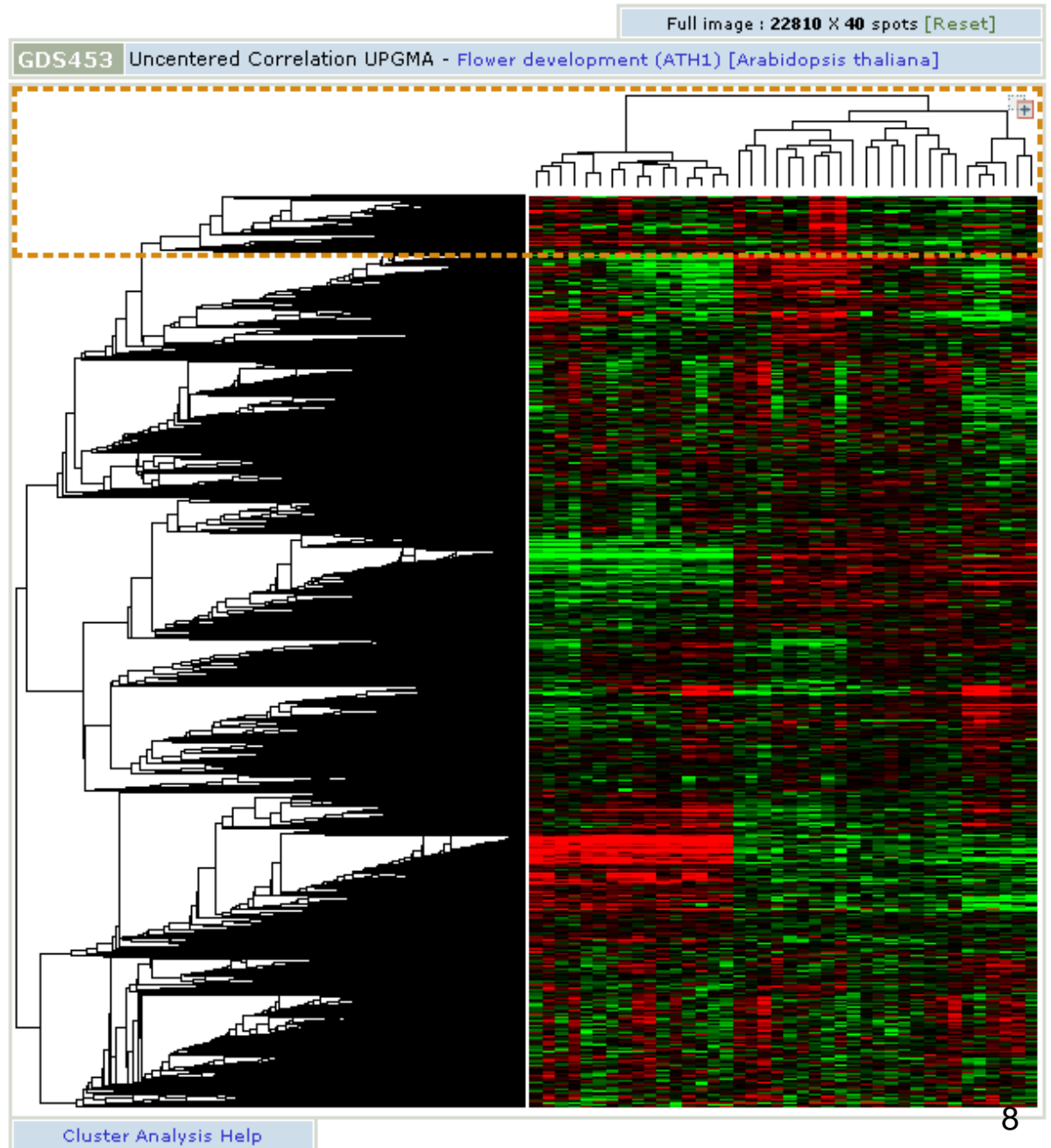


# Картирование транскриптов



# Экспрессия генов – 1. Развитие цветка резуховидки Таля

двойная  
кластеризация  
– на генах и на  
условиях



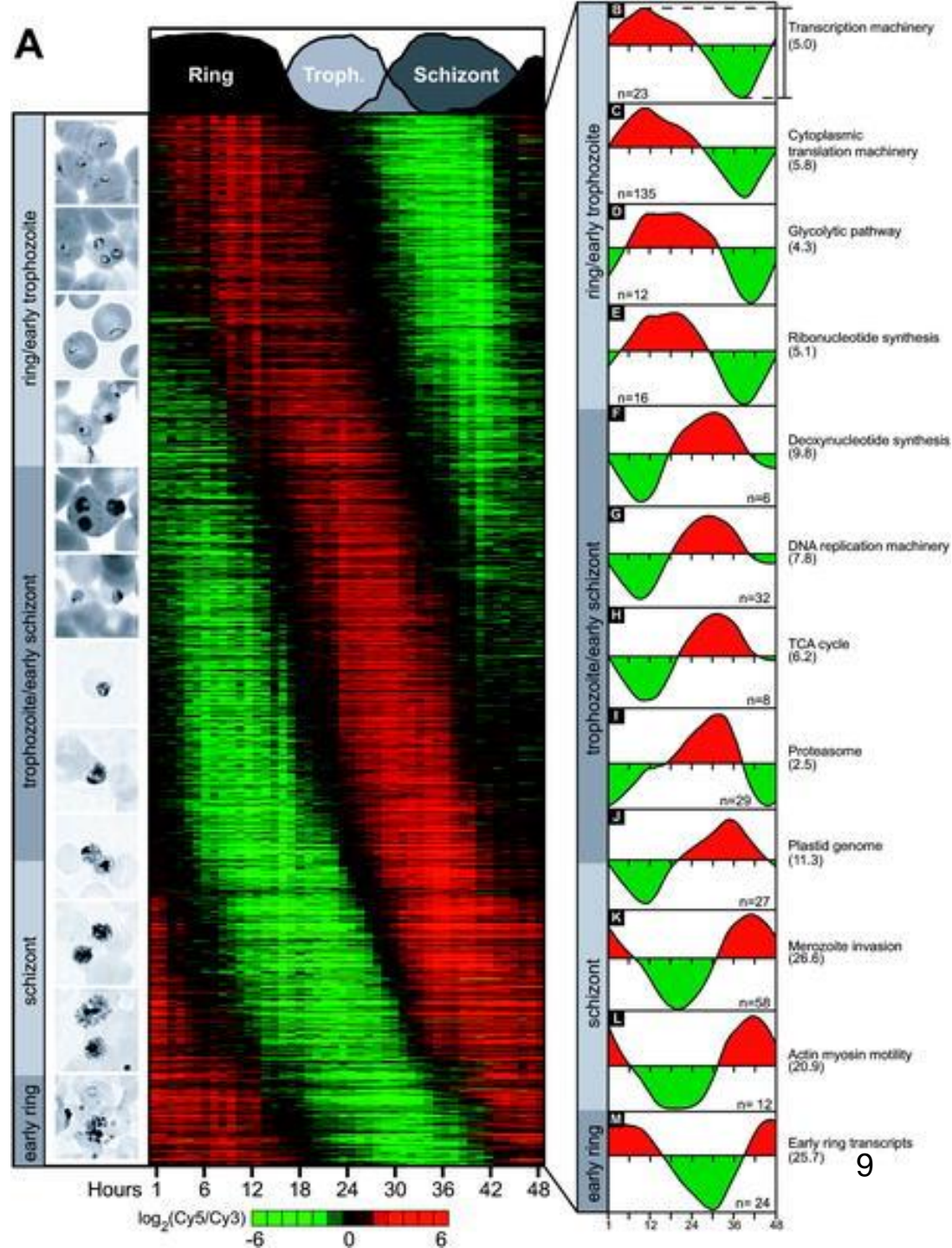


# Экспрессия генов – 2. Цикл развития малярийного плазмодия



The Transcriptome of the Intraerythrocytic Developmental Cycle of *Plasmodium falciparum*

Zbynek Bozdech<sup>1</sup>, Manuel Llinás<sup>1</sup>, Brian Lee Pulliam<sup>1</sup>, Edith D. Wong<sup>1</sup>, Jingchun Zhu<sup>2</sup>, Joseph L. DeRisi<sup>1</sup>



# Транскриптомы: типичные задачи

- Когорты пациентов
- Опухоль – ткань
  
- Дифференциальный диагноз
- Прогноз
- Подбор лечения
  
- Деконволюция

# Single cell геномика и транскриптомика

- Геномика: история клеток
  - Нейроны
  - Рак
- Транскриптомика: различия клеток
  - Intrinsic noise
  - Импринтинг
  - Эмбриология
  - Рак: клоны

# Проблемы

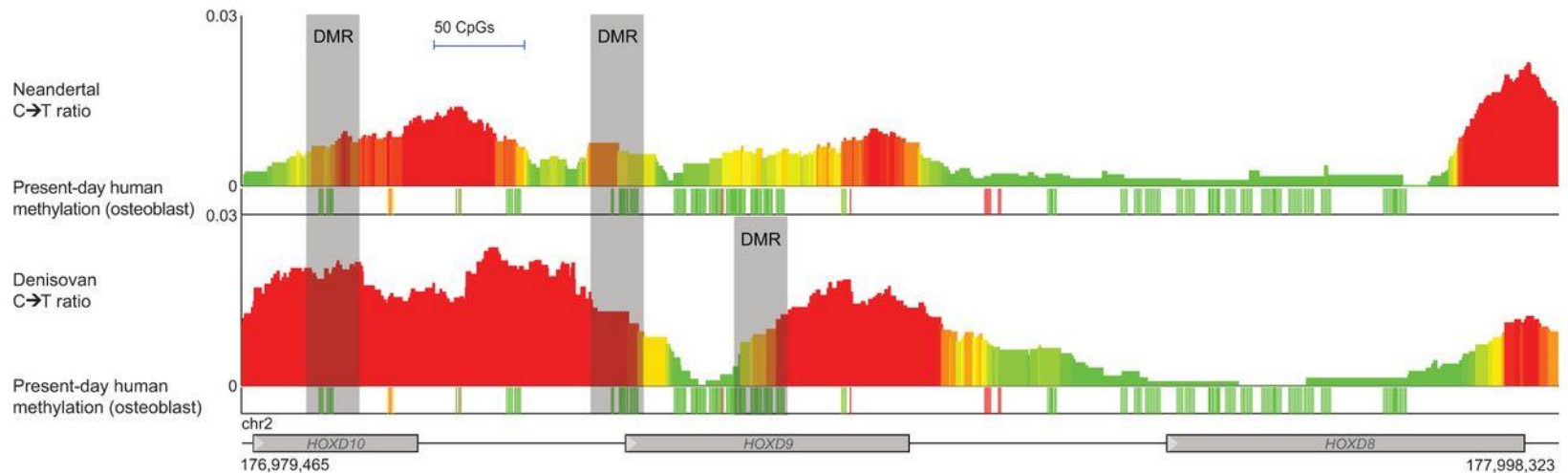
- Шум: каждое отдельное наблюдение ненадежно
- Множественное тестирование. FDR
- Большая значимость – слабый эффект
- Взаимозависимость параметров
- Абстракция: работаем с конструктами, а не с реальными объектами

- Предсказание (медицина)

vs

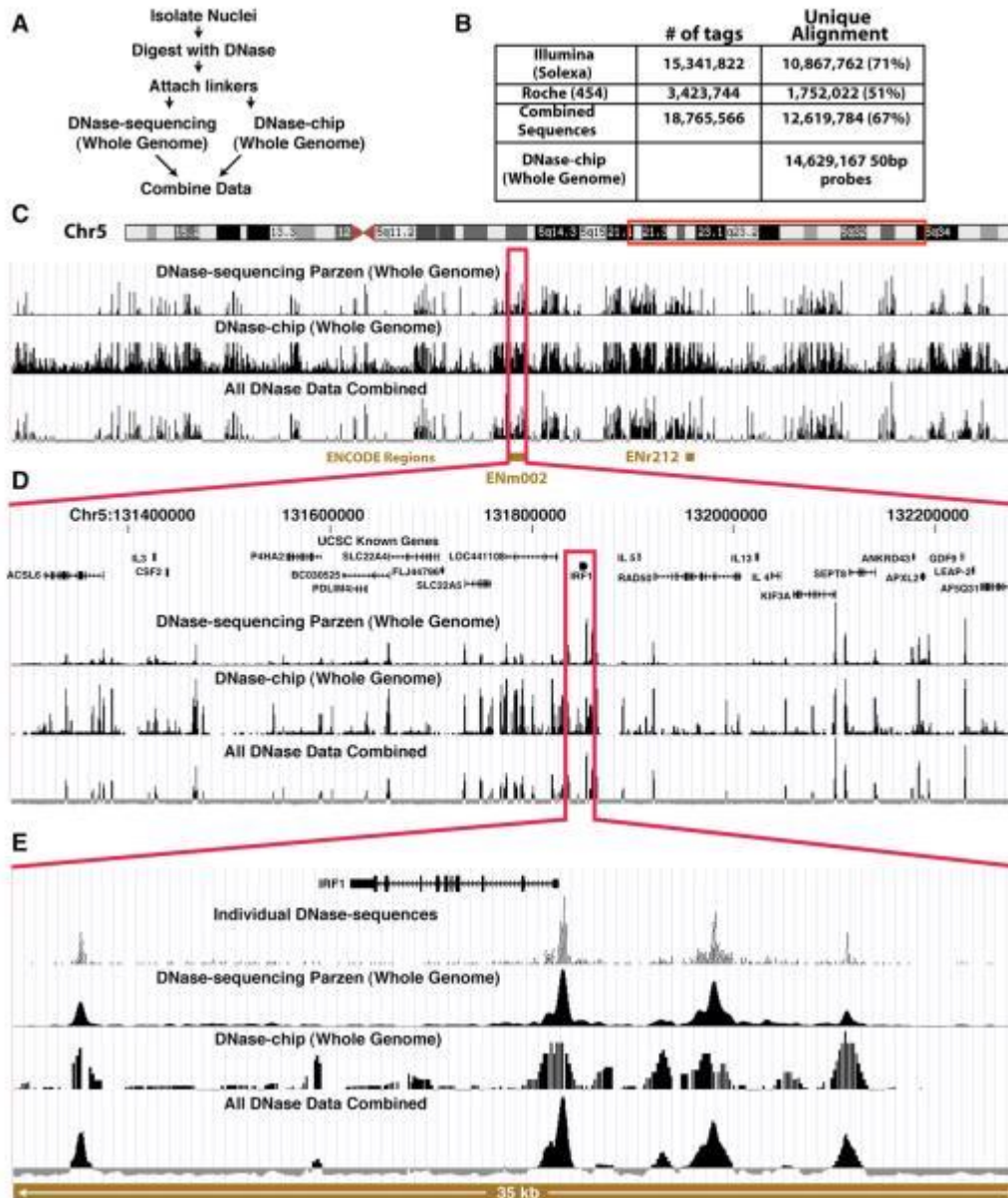
- Feature extraction (биология)

# Метилирование

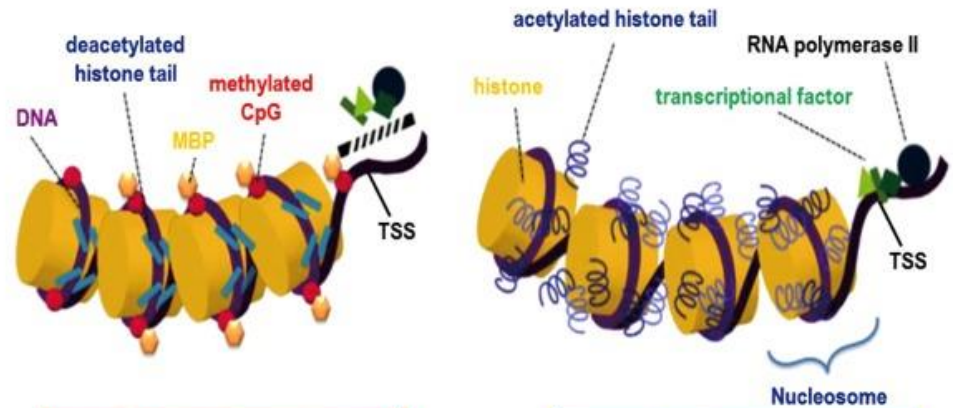
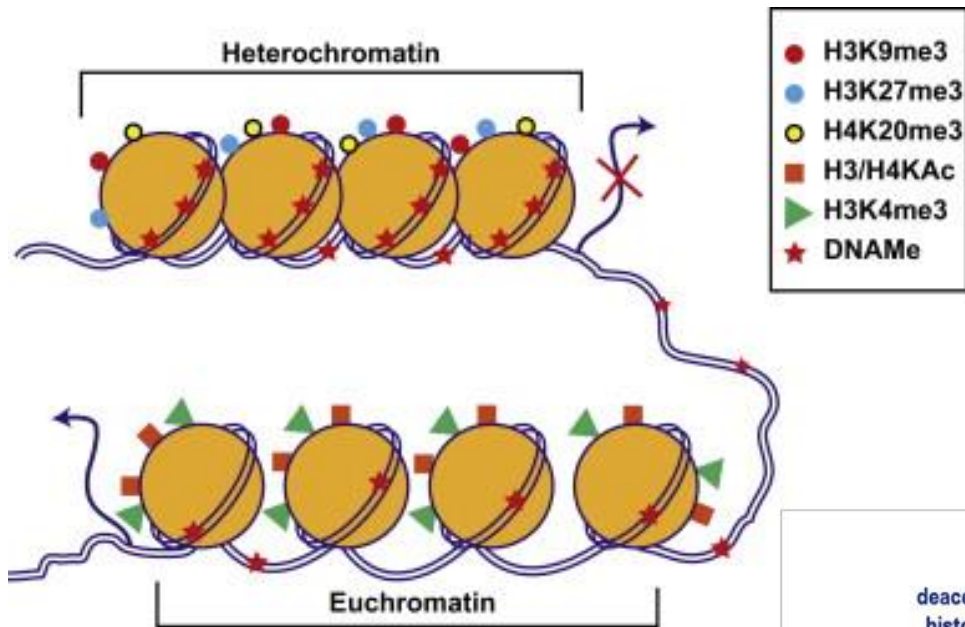


The HOXD cluster is hypermethylated in archaic humans. C→T ratio in archaic humans, and present-day human methylation

# DNase I – открытый хроматин



# Модификации гистонов

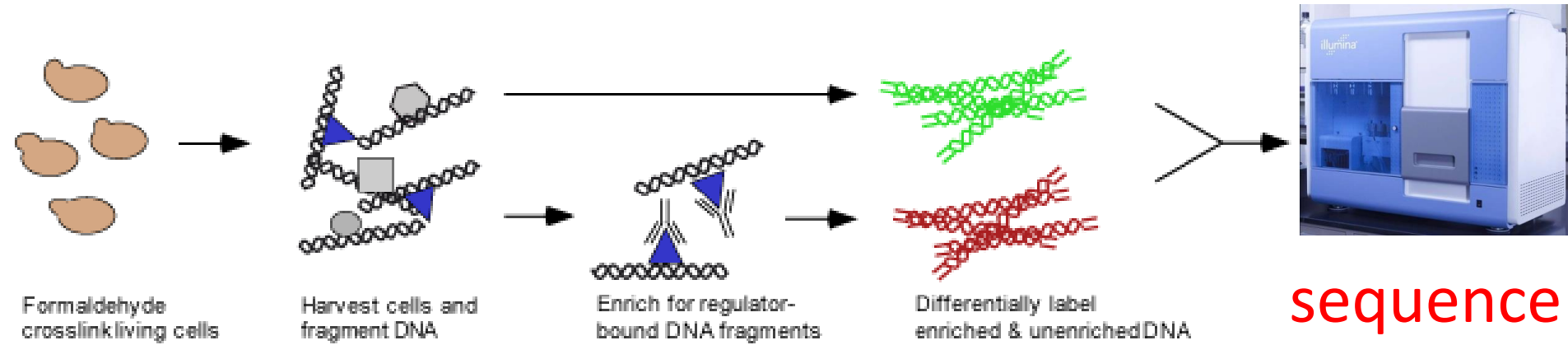


**Heterochromatin**  
closed chromatin  
conformation  
: repression

**Euchromatin**  
open chromatin  
conformation  
: activation

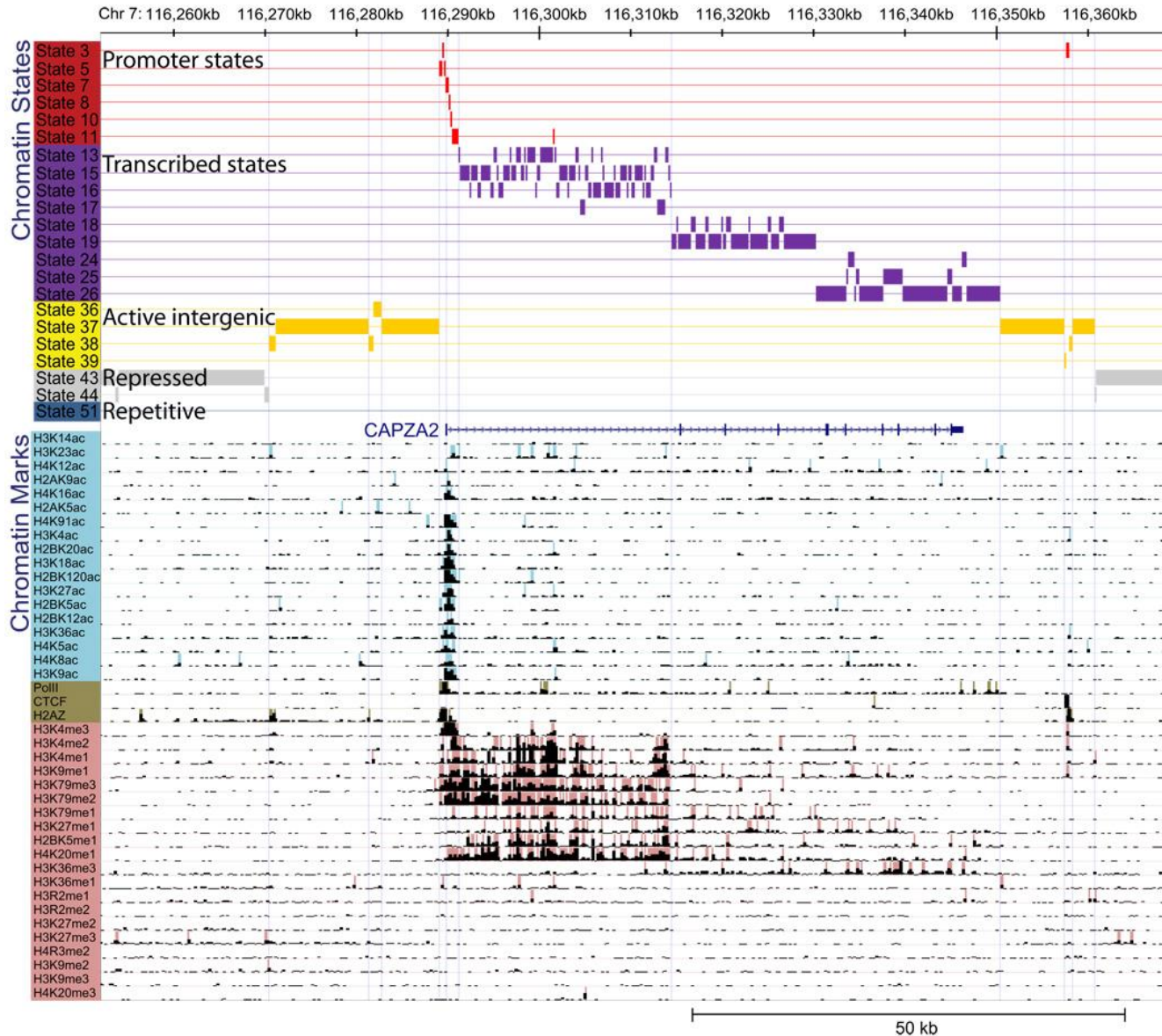
TRANSCRIPTION

# ChIP-Seq

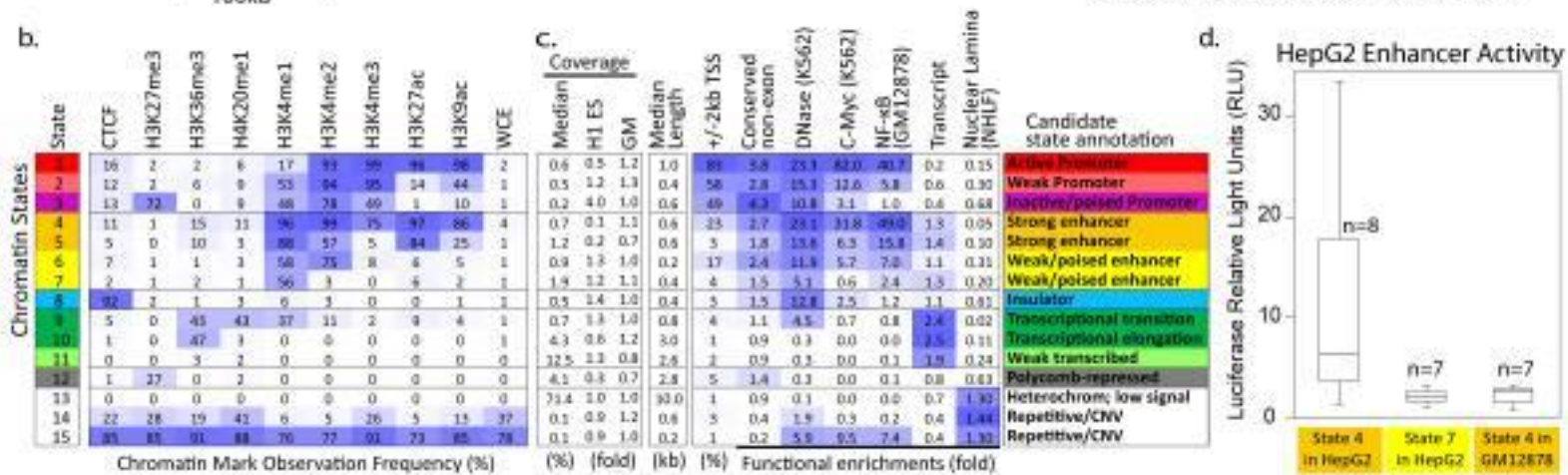
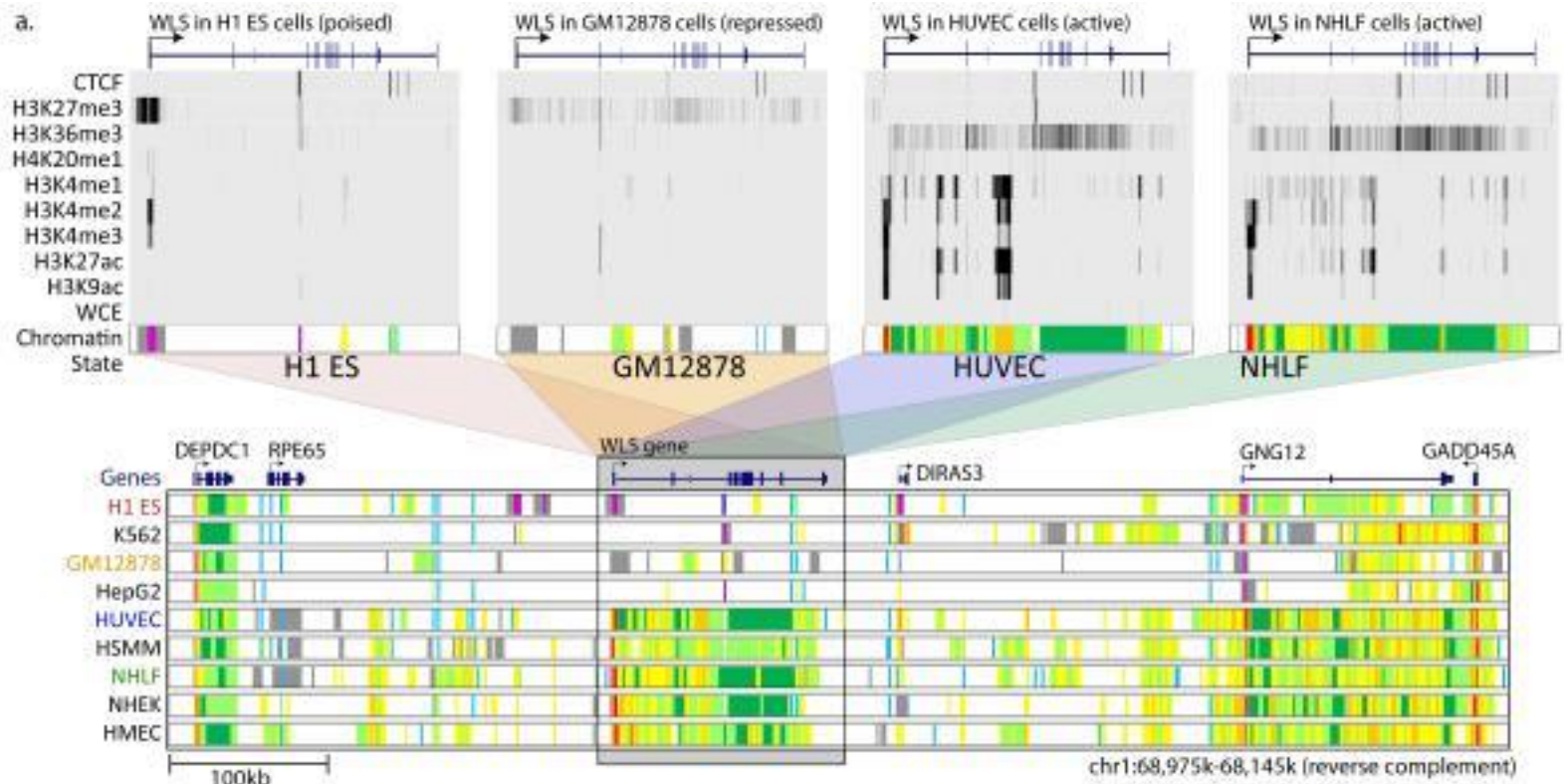




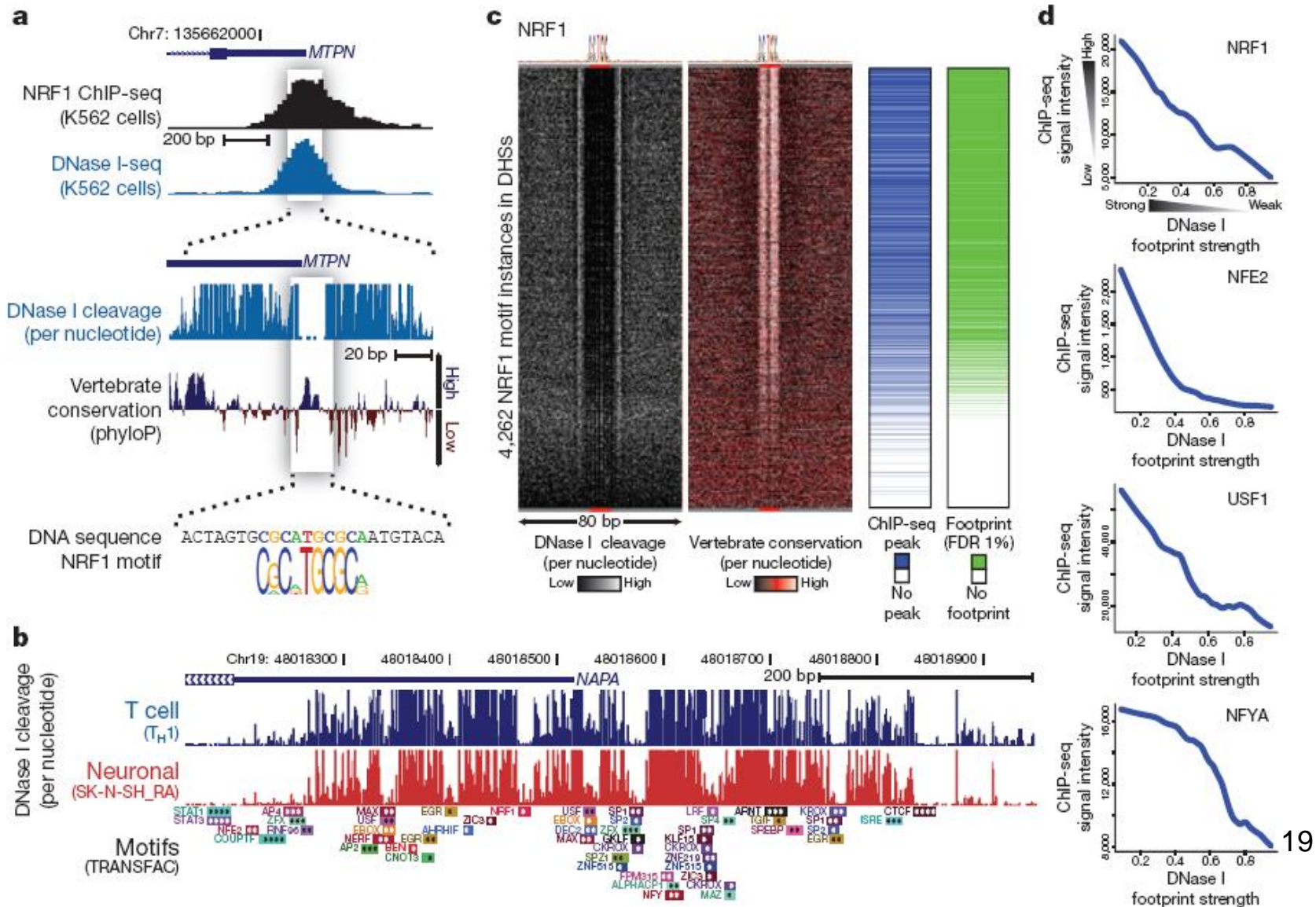
# Индивидуальные маркеры и интегрированные состояния



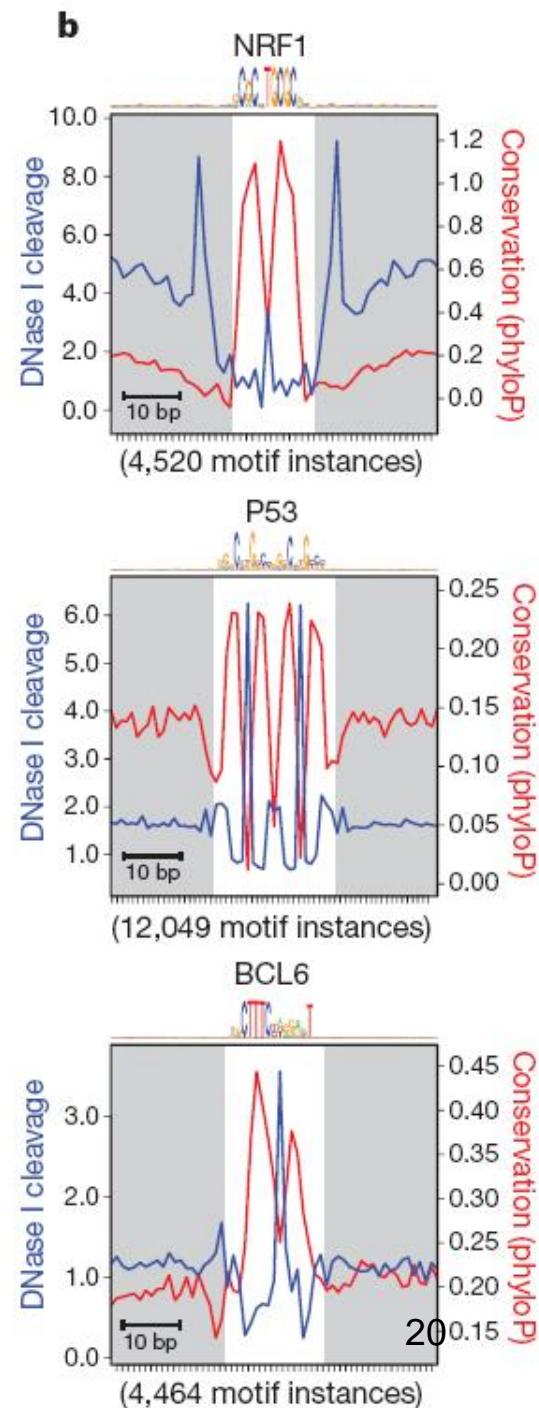
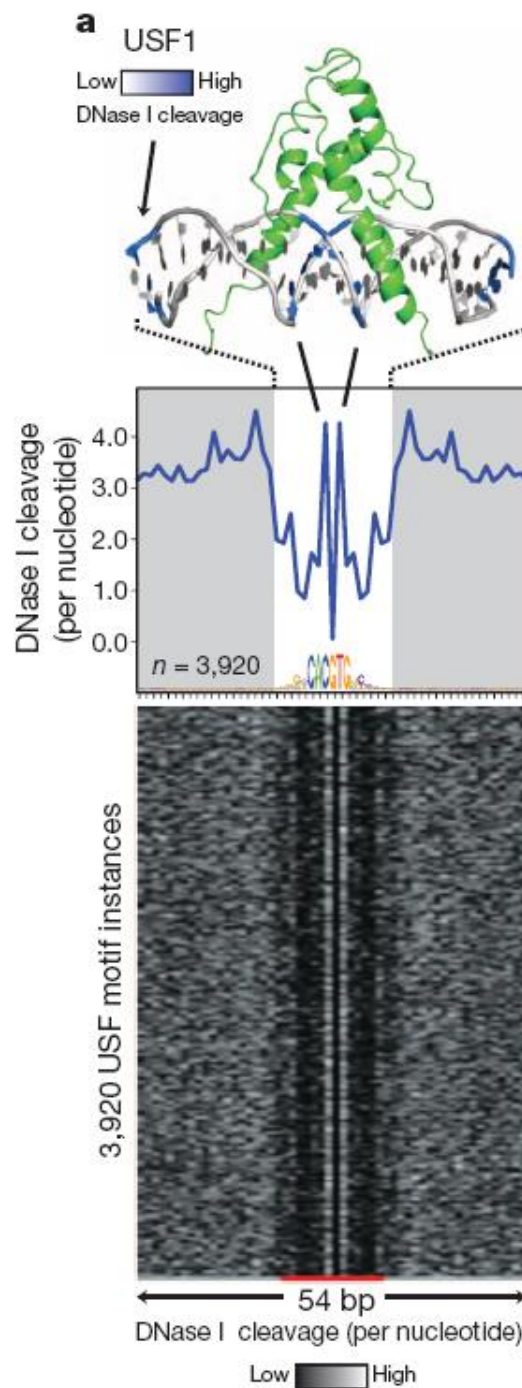
# Другой пример: сравнение клеточных линий



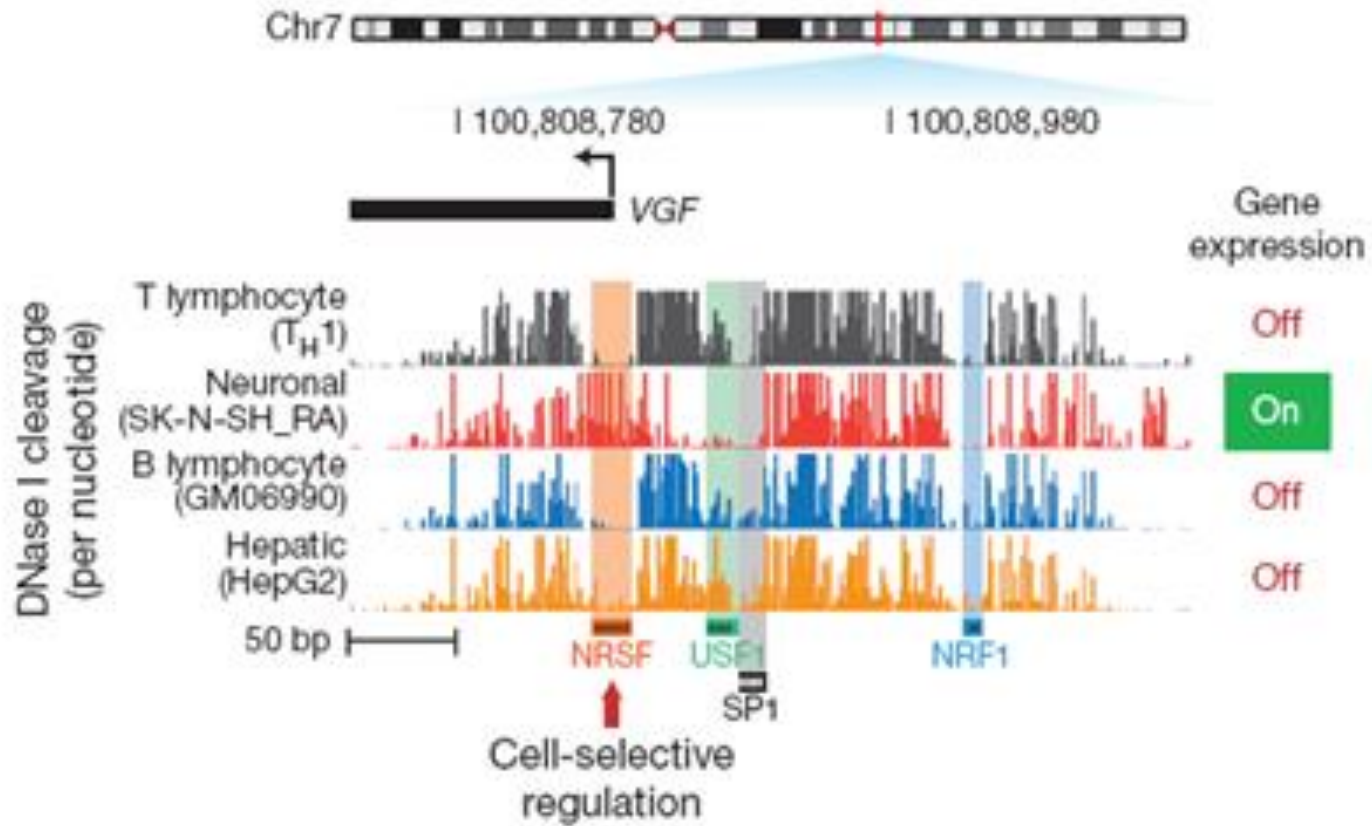
# Доступность хроматина для ДНКазы I



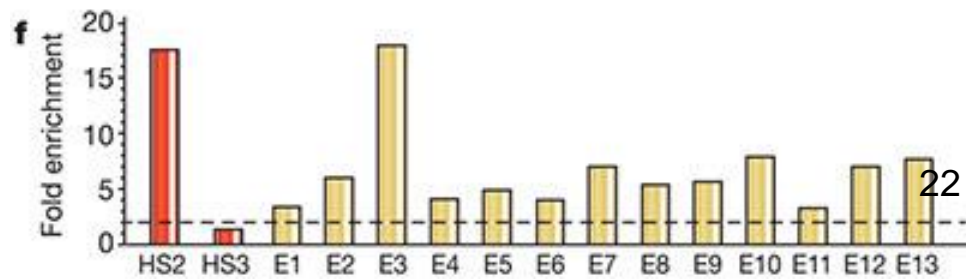
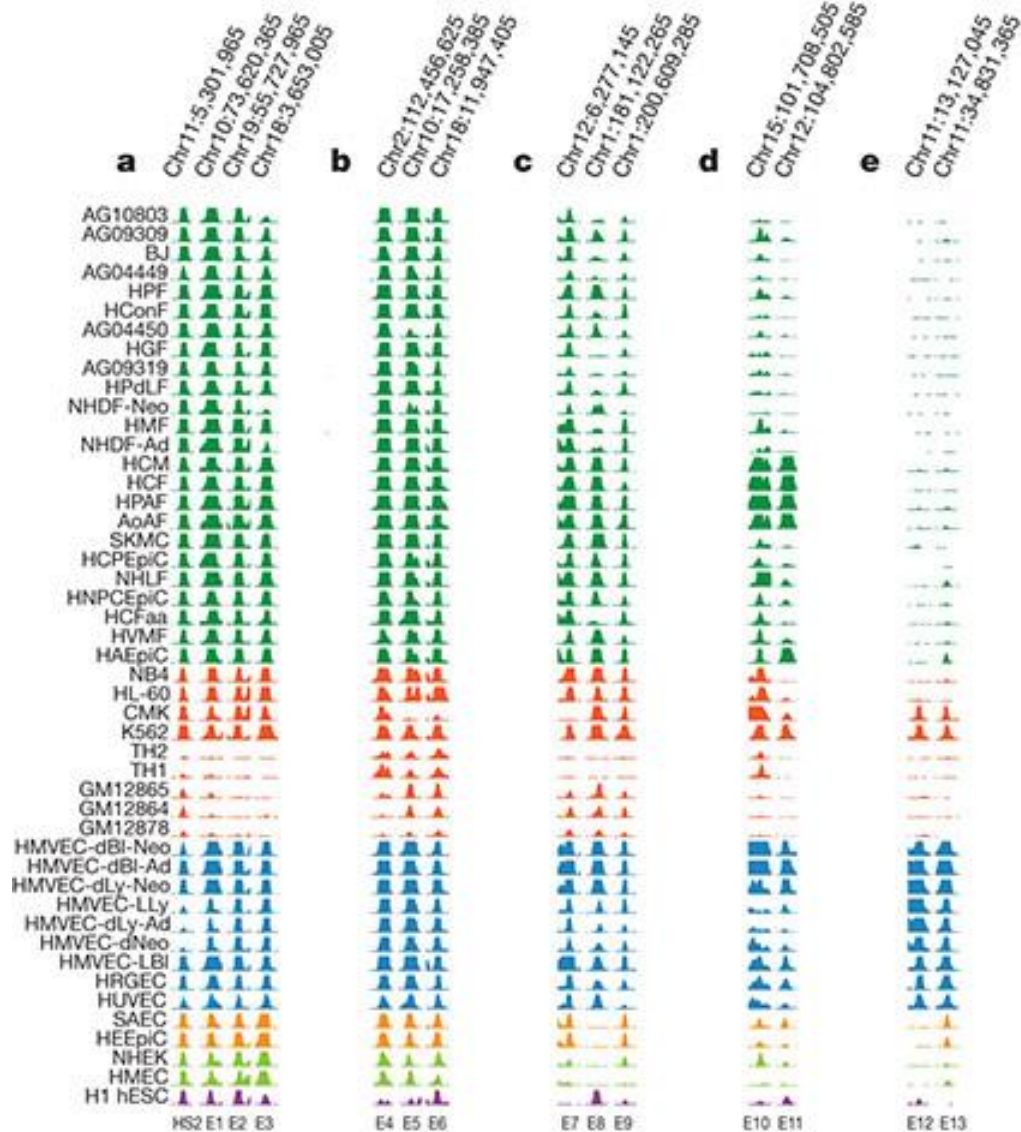
# Доступность хроматина для ДНКазы I и структура комплекса фактора транскрипции - ДНК



# Тканевая специфичность

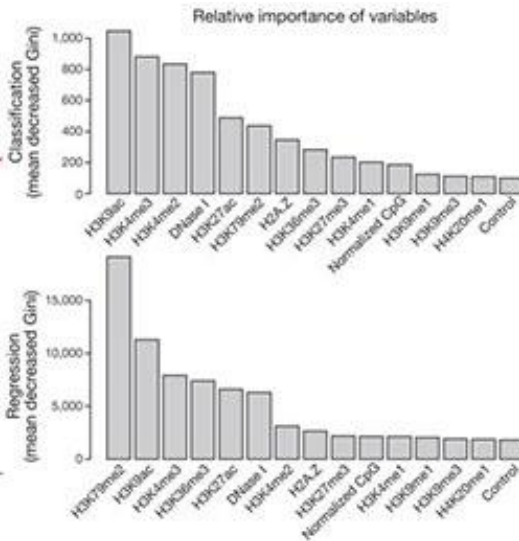
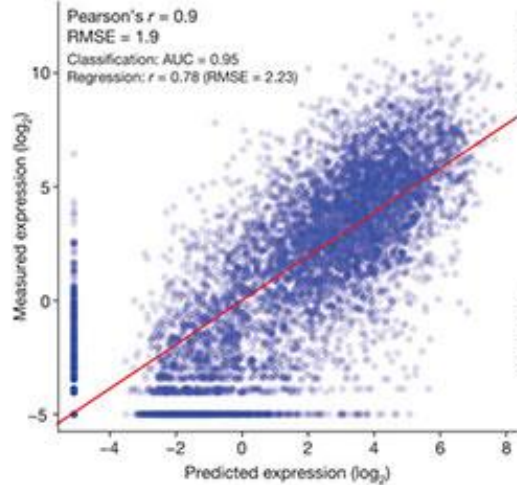


# Кластеризация энхансеров по открытости в разных клеточных линиях

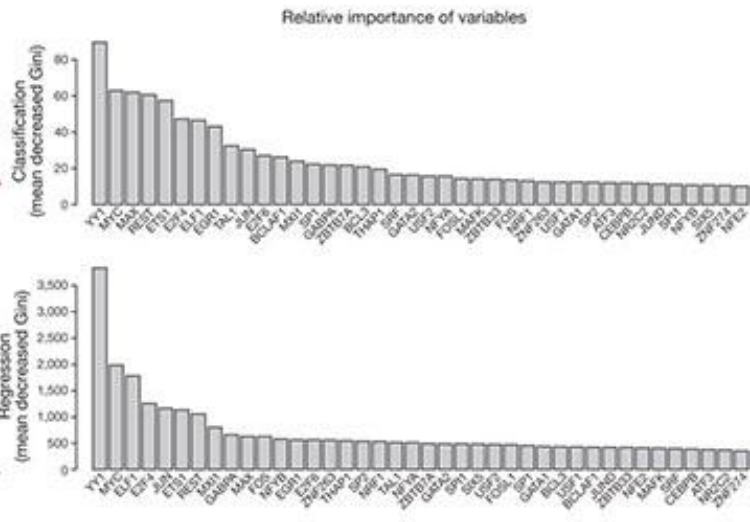
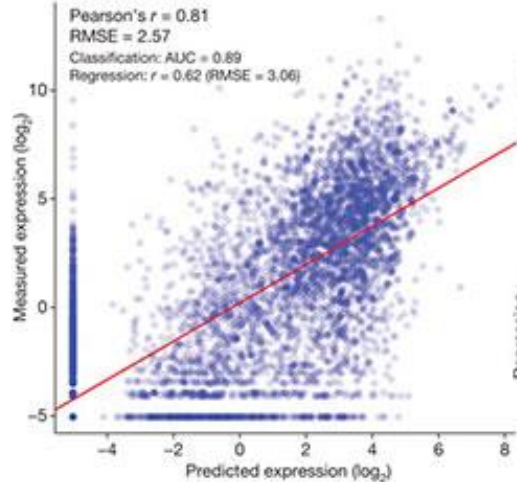


# Предсказание уровня экспрессии по модификациям гистонов и по связыванию факторов транскрипции

**a** CAGE poly(A)<sup>+</sup> K562 whole cell

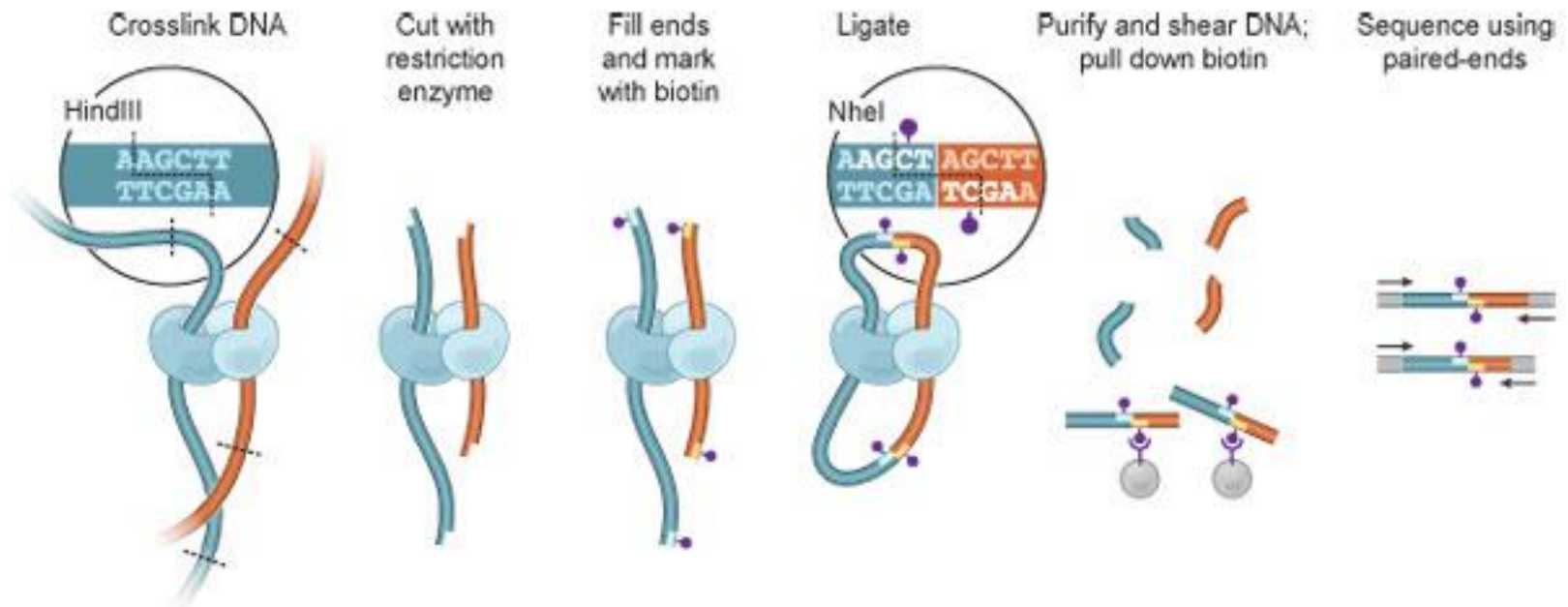


**b** CAGE poly(A)<sup>+</sup> K562 whole cell



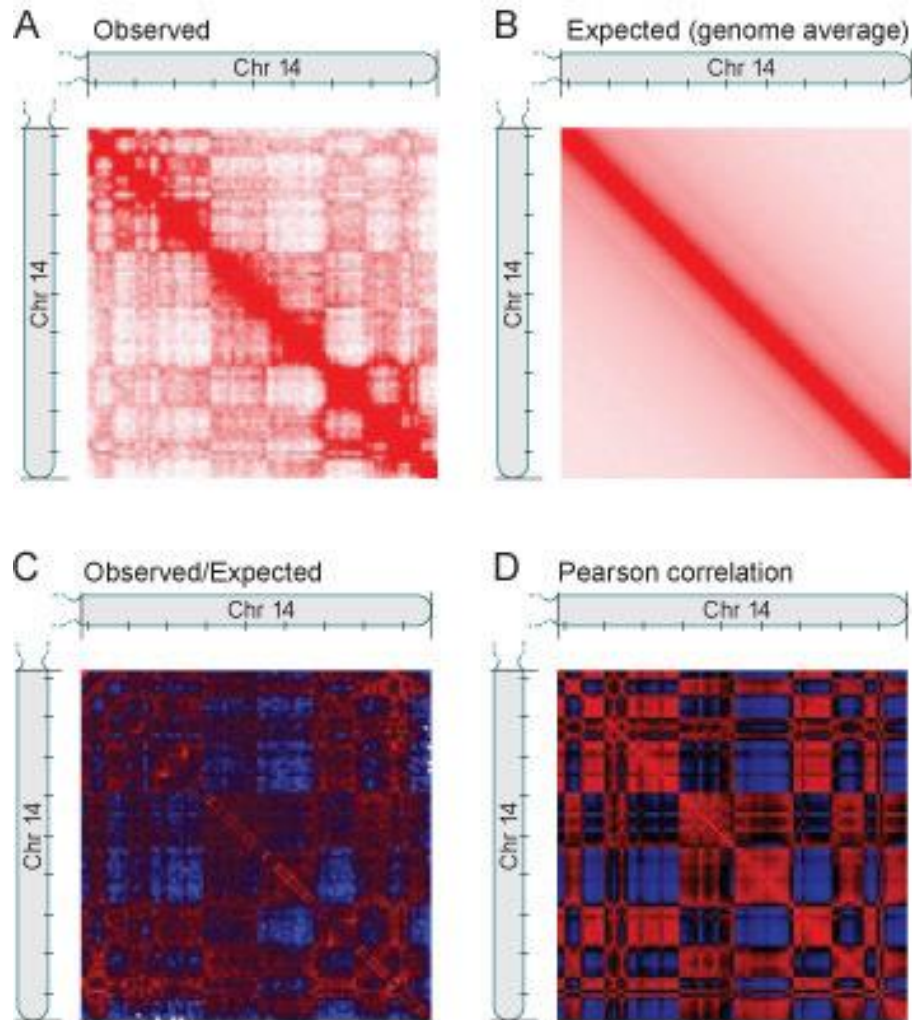
# HiC

E.Lieberman-Aiden ... J.Dekker, Science, 2009





# Карты HiC

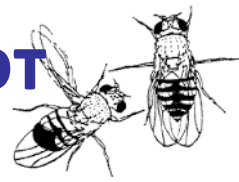


NB: это не геометрическая близость в 3D, а частота контактов.

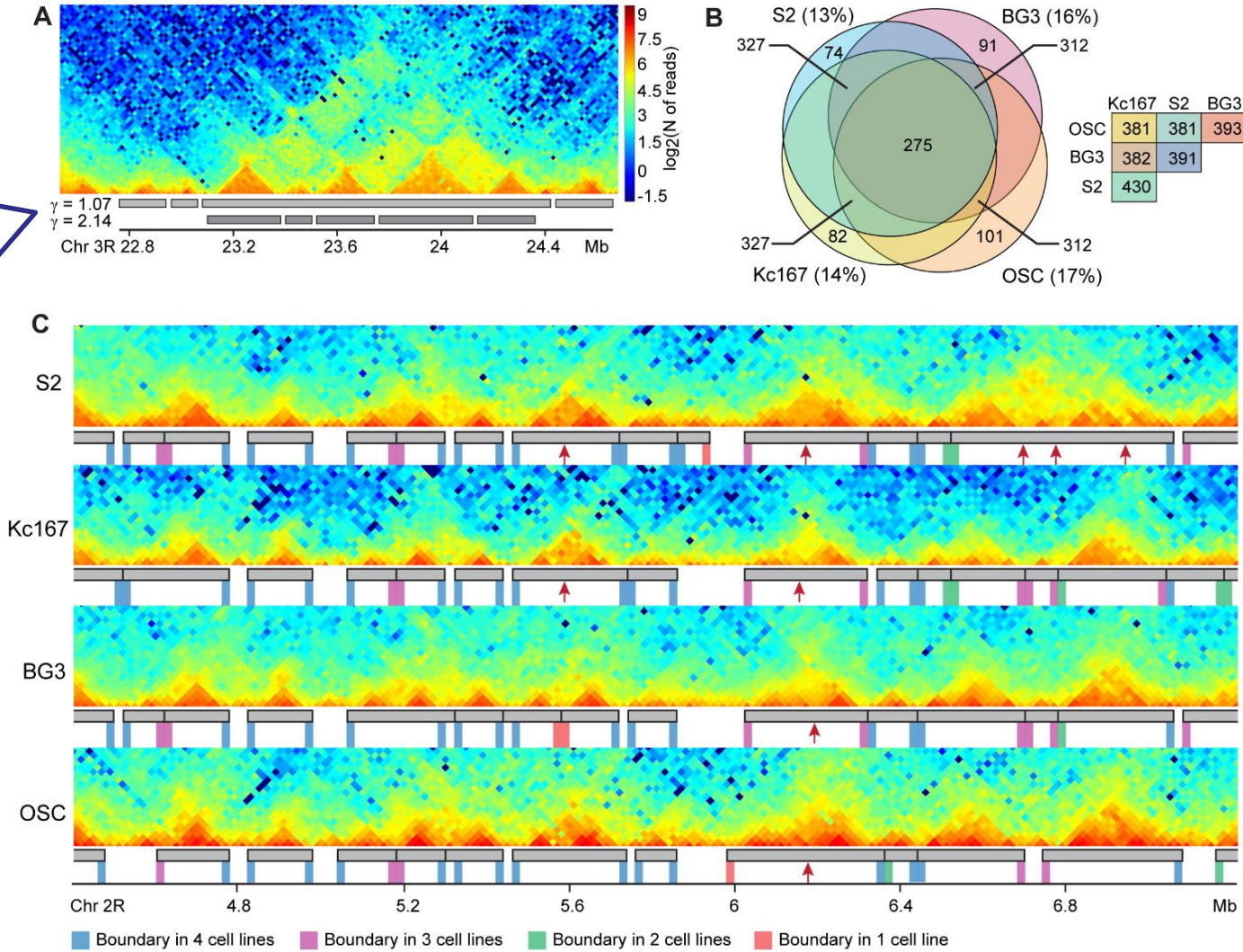
Корреляция (общие соседи) как прокси для геометрии

Dekker's group

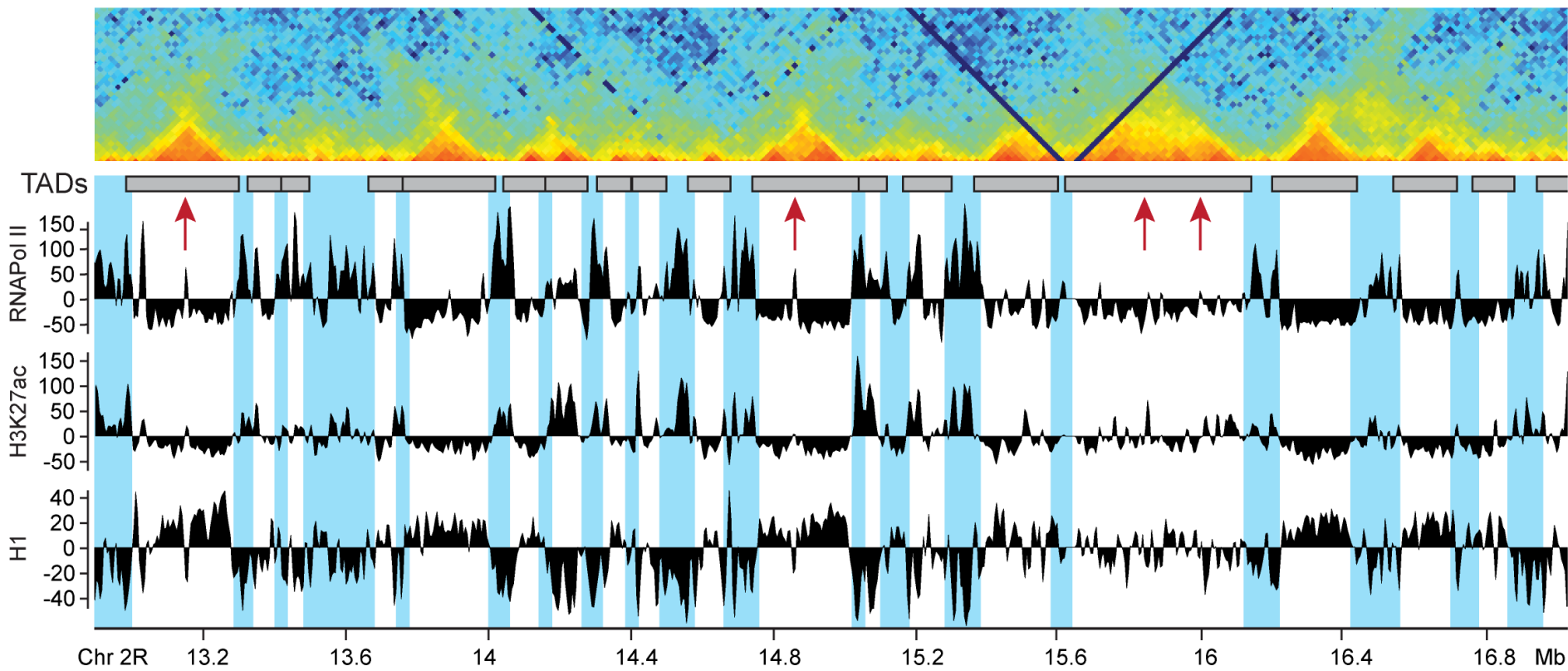
# Границы ТАДов в основном совпадают между клеточными линиями



Definition of TADs depends on a parameter ( $\gamma$ ).  
 Low  $\gamma \Rightarrow$  larger TADs

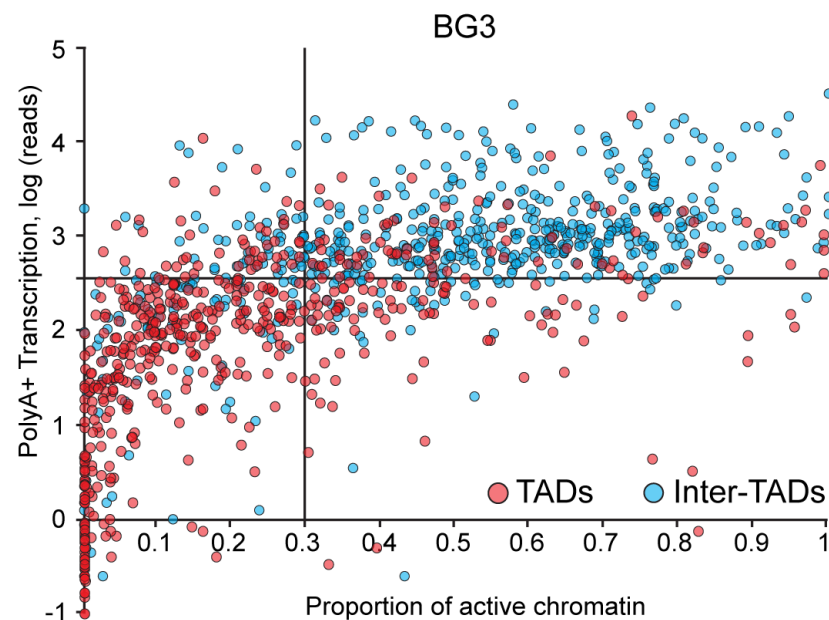
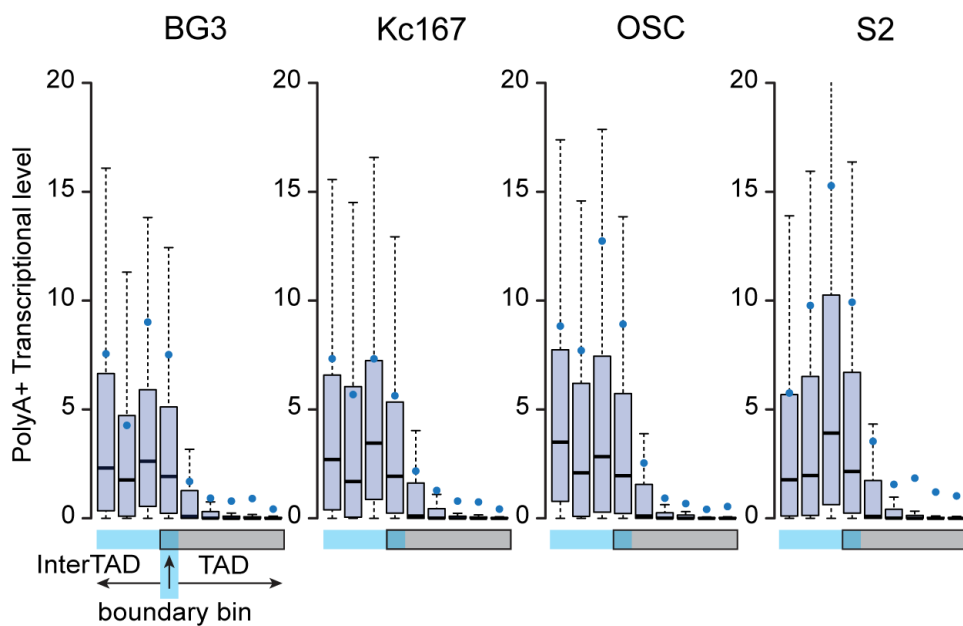


# Меж-ТАДы обогащены активными хроматиновыми метками и *pol II* и обеднены репрессивными метками



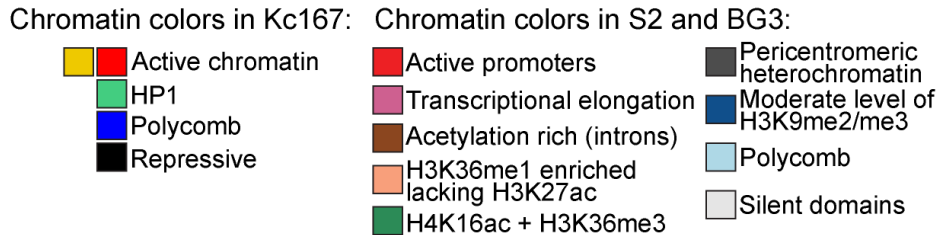
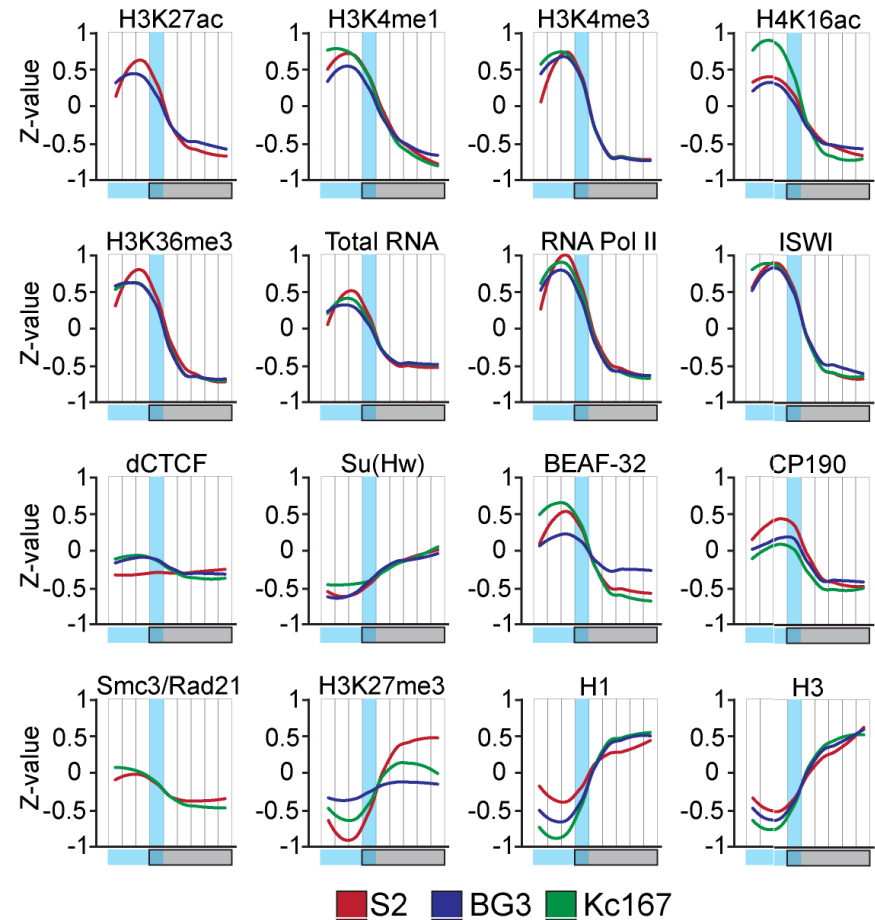
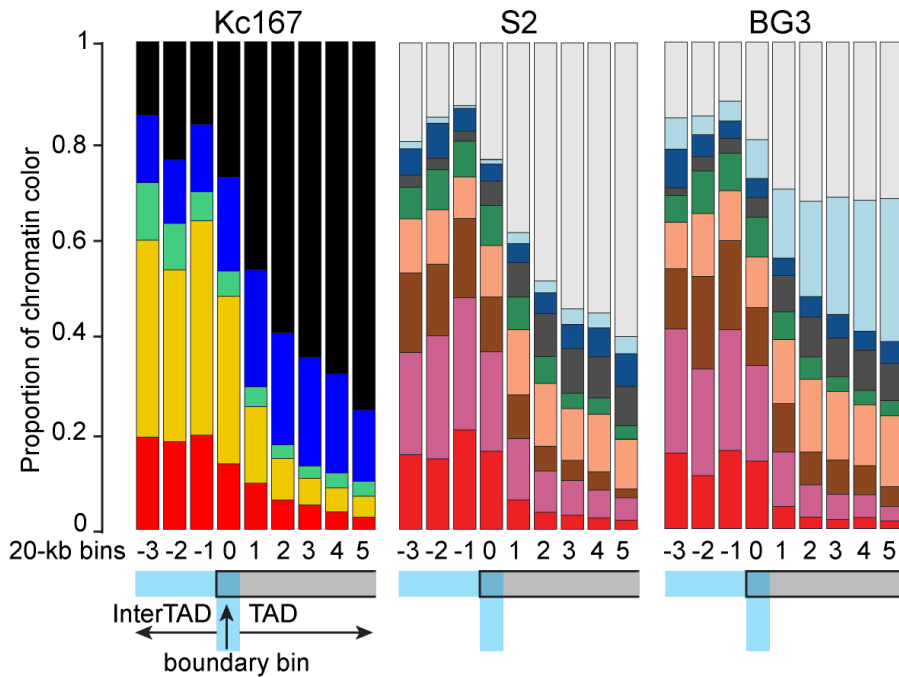


# Уровень транскрипции и доля активного хроматина в ТАДах и меж-ТАДах



Log10 normalized read counts

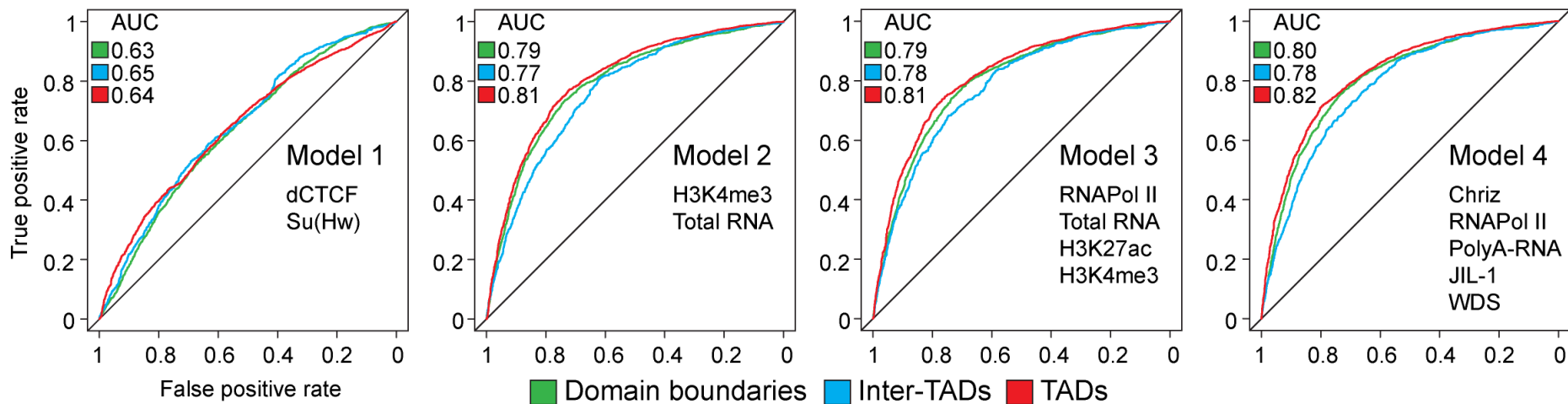
# Распределение хроматиновых меток на границах ТАДов





# Предсказание границ по меткам и сайтам связывания инсуляторов (попытка feature extraction)

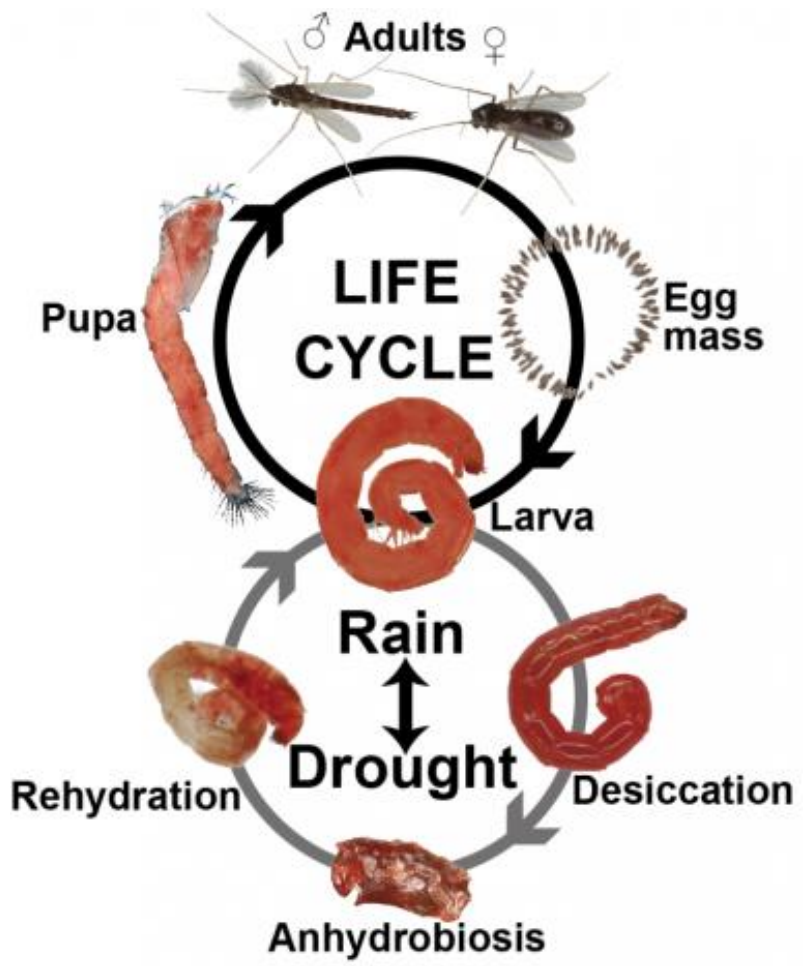
S2 ROC-curves



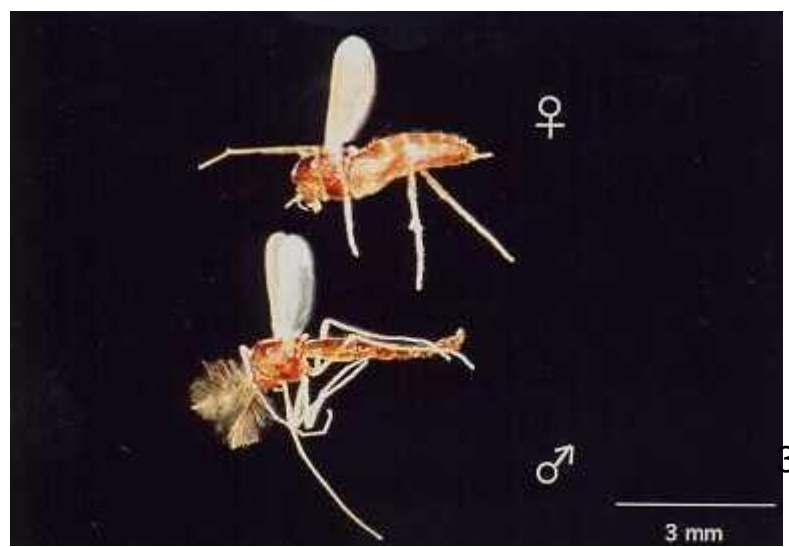
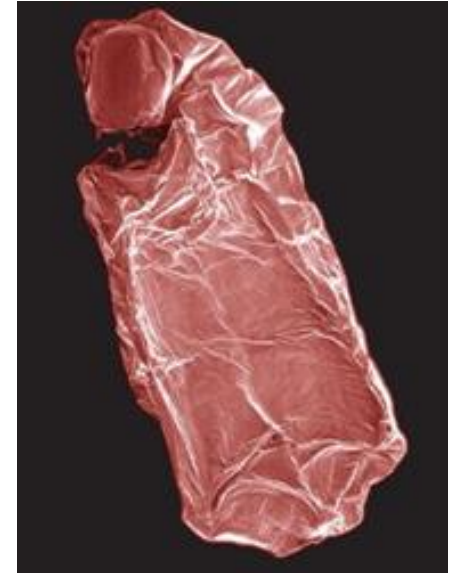
## A Fly Larva that tolerates Dehydration and Temperatures of $-270^{\circ}$ to $+102^{\circ}$ C.

H. E. HINTON

Department of Zoology, University of Bristol.



Cornette & Kikawada, 2011, IUBMB Life



- Жизненный цикл

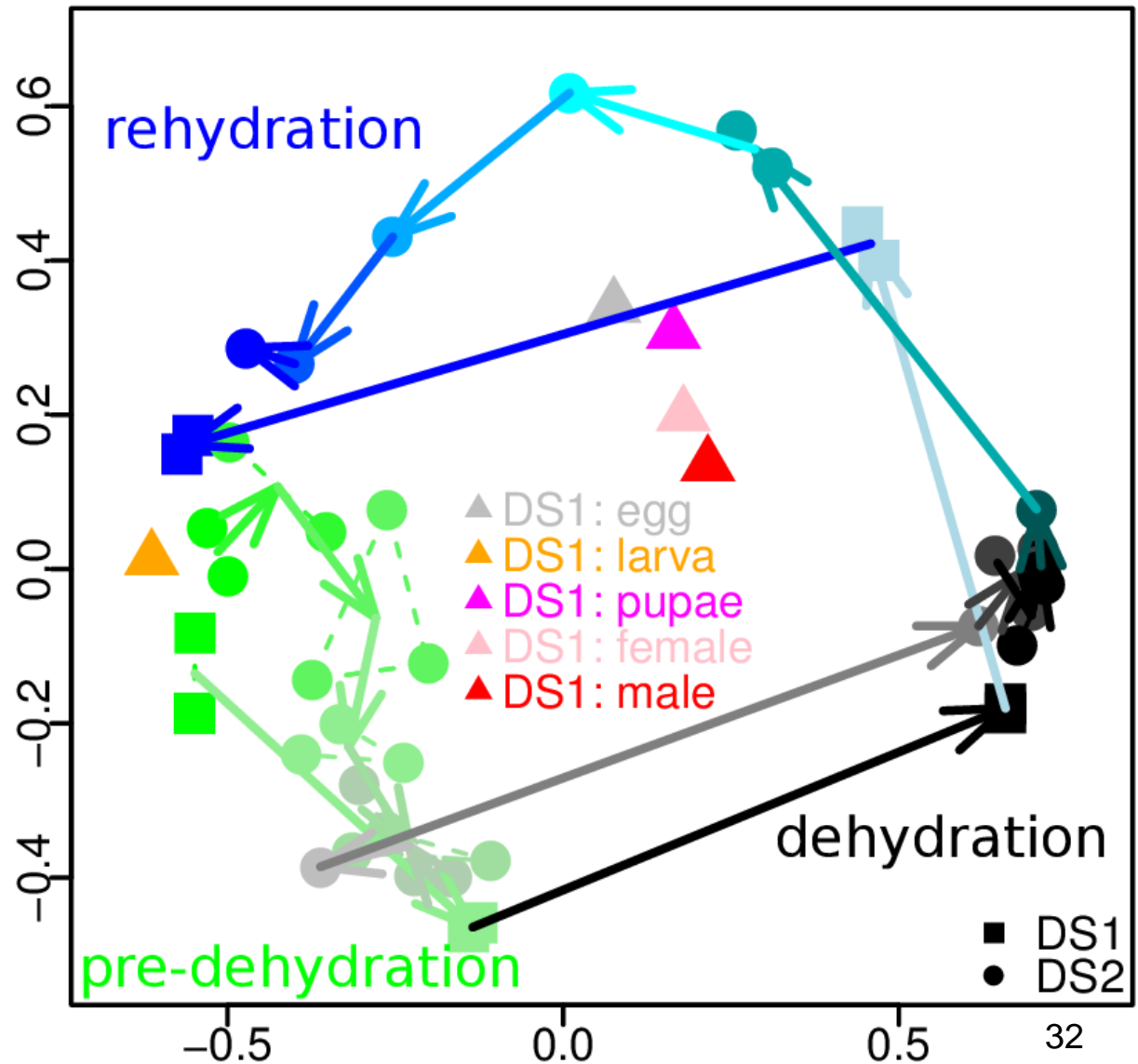
- яйцо
- личинка
- куколка
- самец и самка

- Цикл

высыхание –  
размачивание

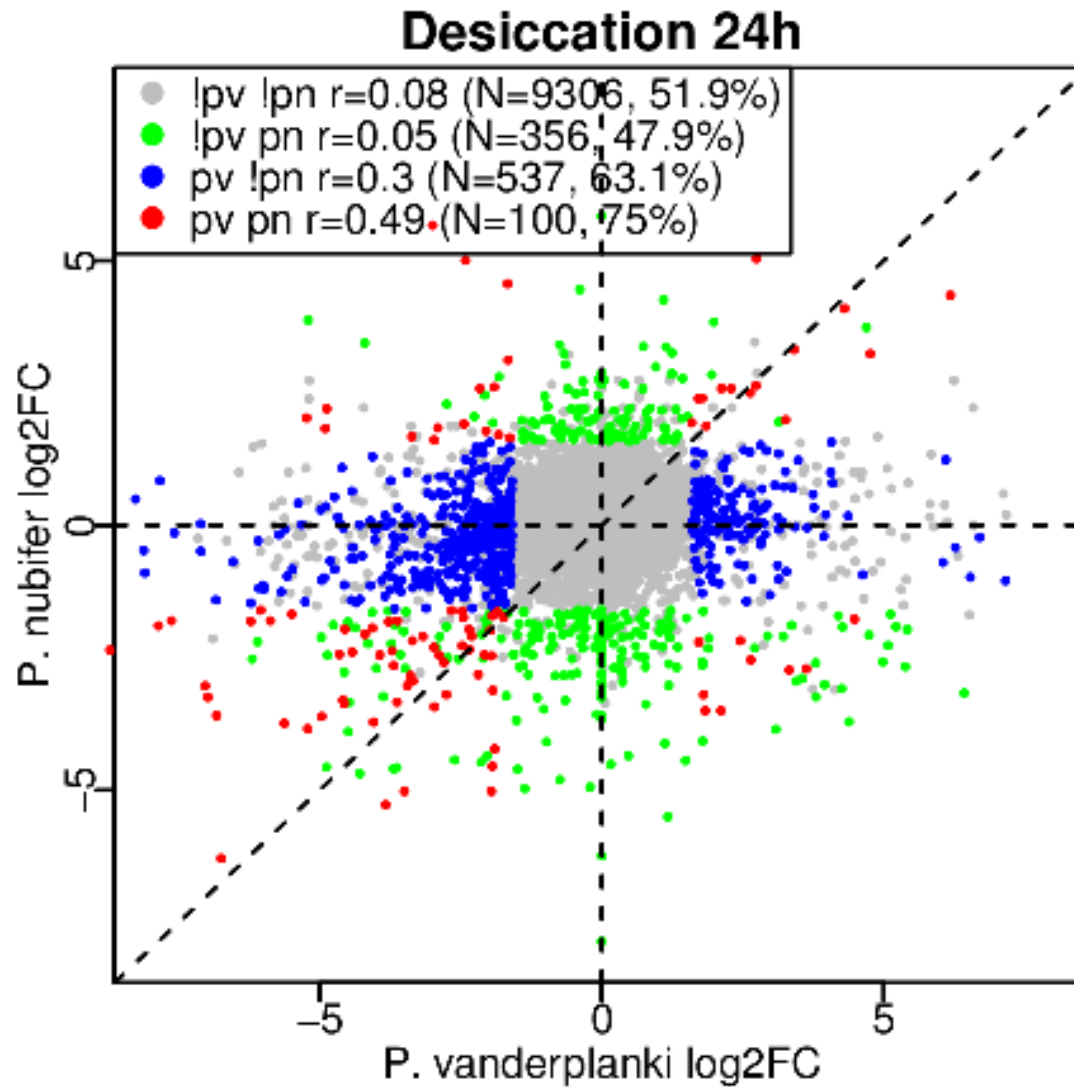
- высыхание:  
0, 24, 48 ч.
- размачивание:  
3, 24 ч.
- *P. nubifer*:  
0 и 24 ч.

## RNA-seq

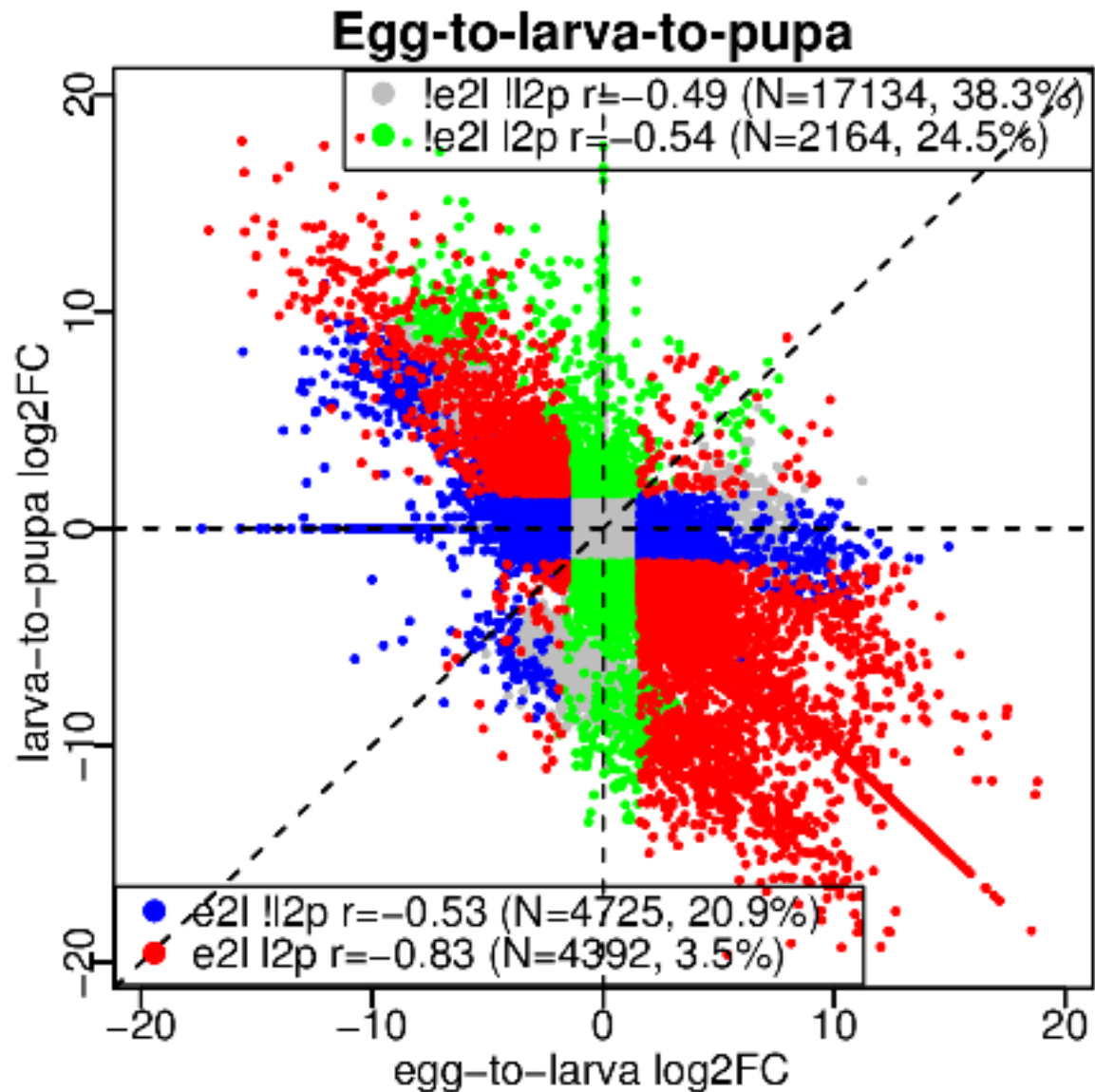




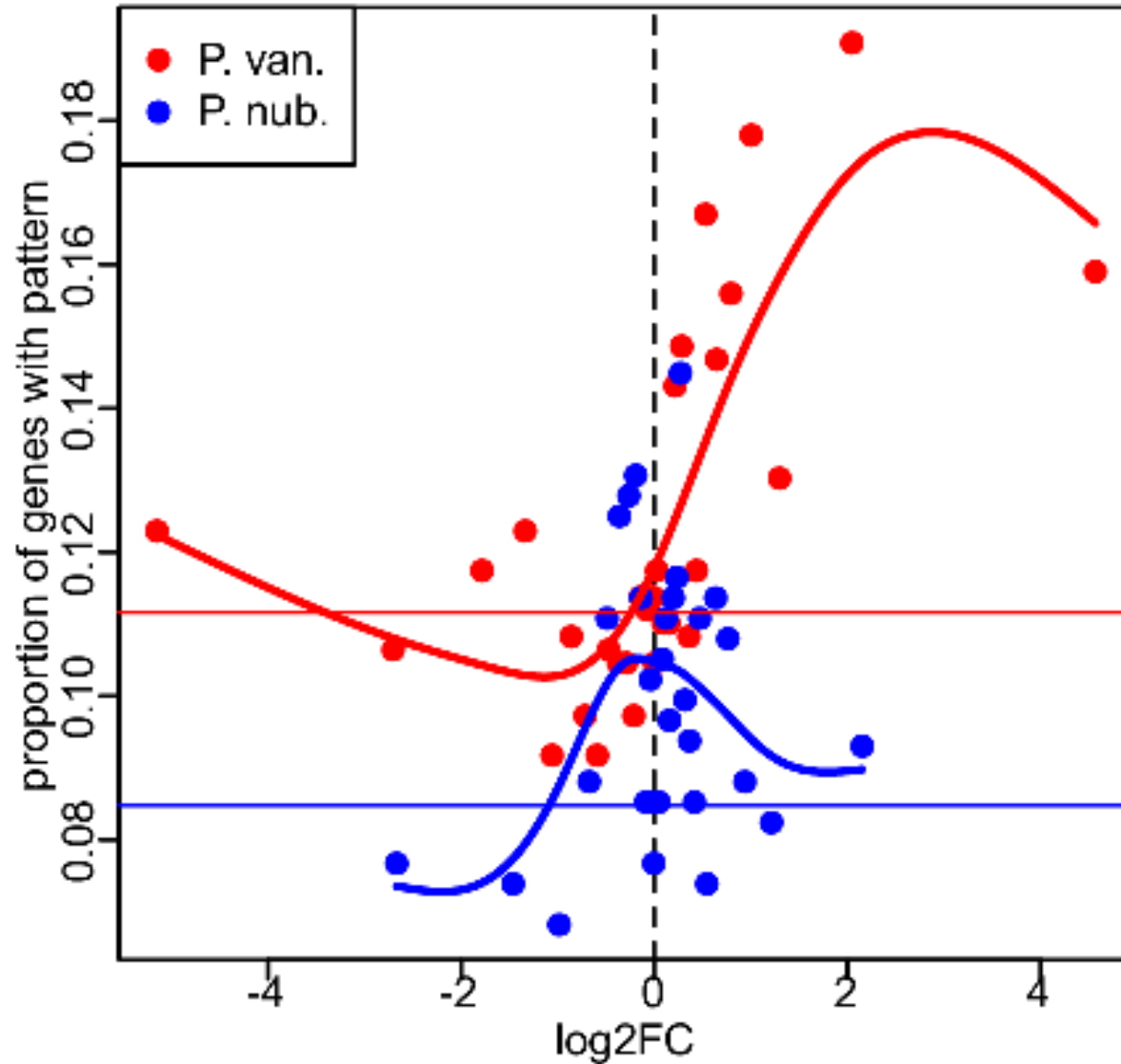
# Высыхание – два комара



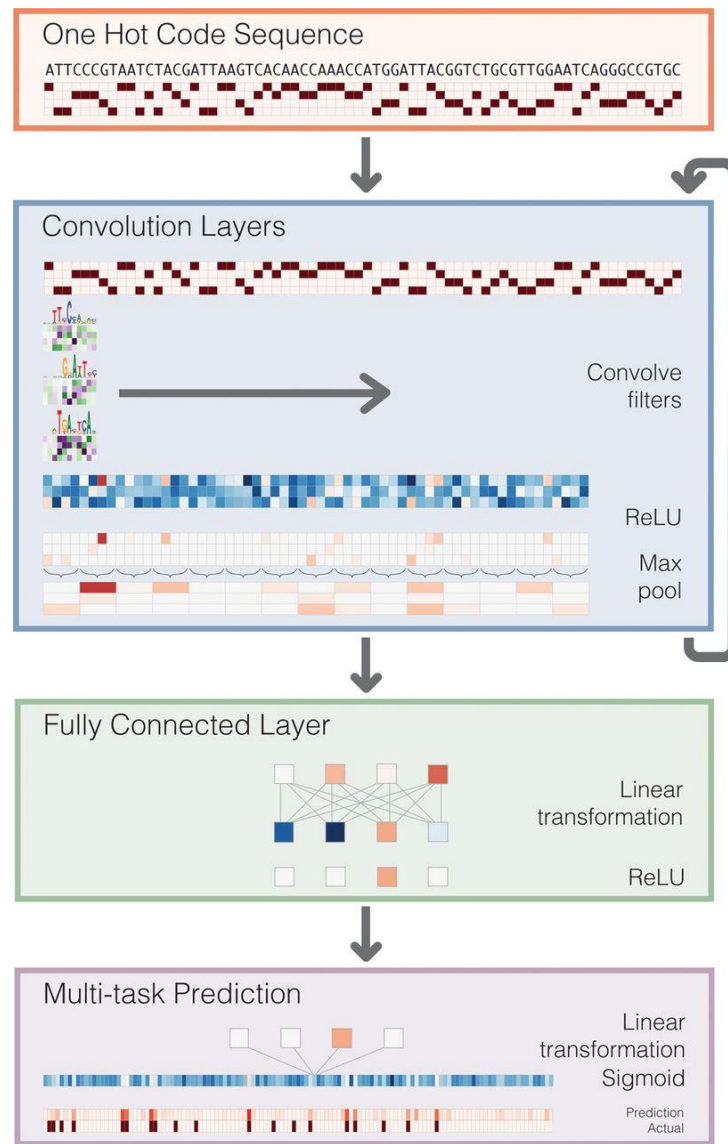
# Развитие: яйцо – личинка – куколка



# TCTAGAA => рост экспрессии при высыхании



# Глубокая сверточная нейронная сеть (CNN)



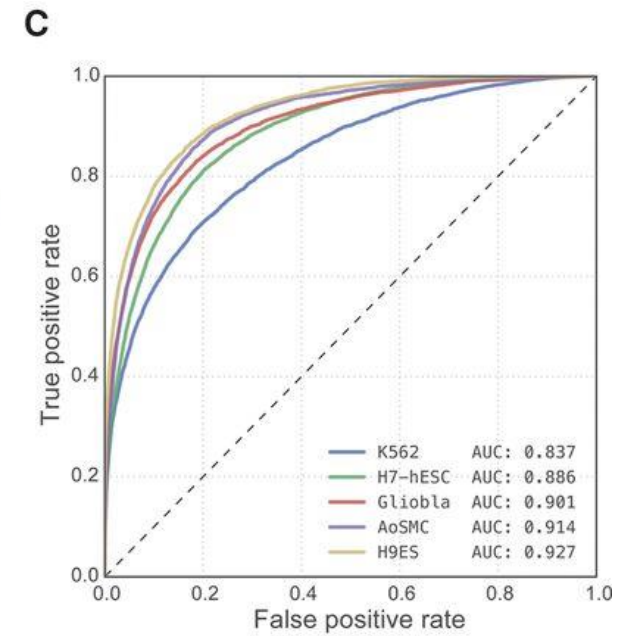
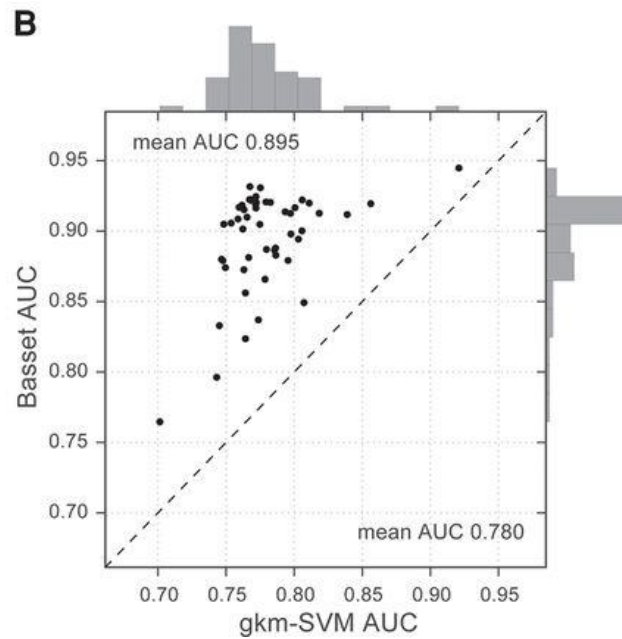
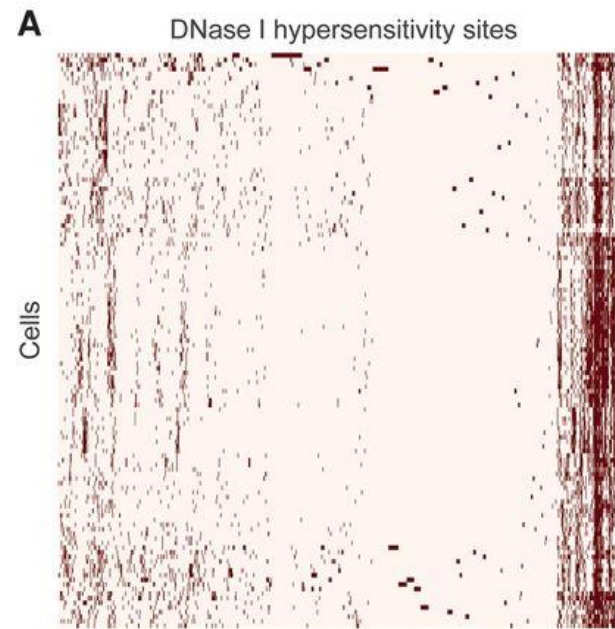
Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks

David R. Kelley et al. *Genome Res.* 2016;26:990-999



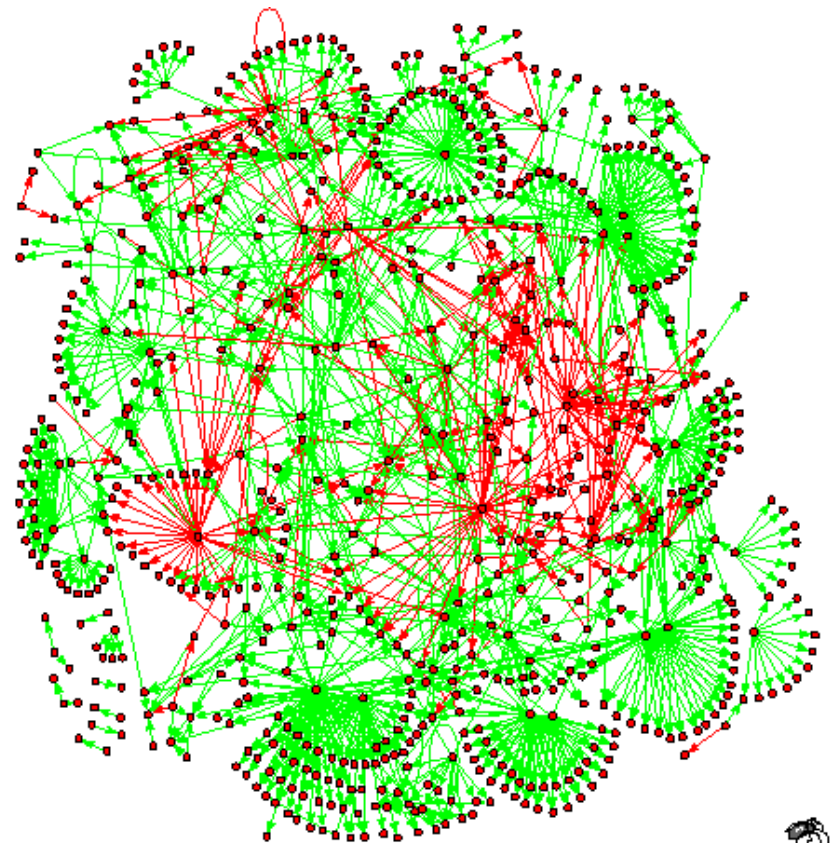
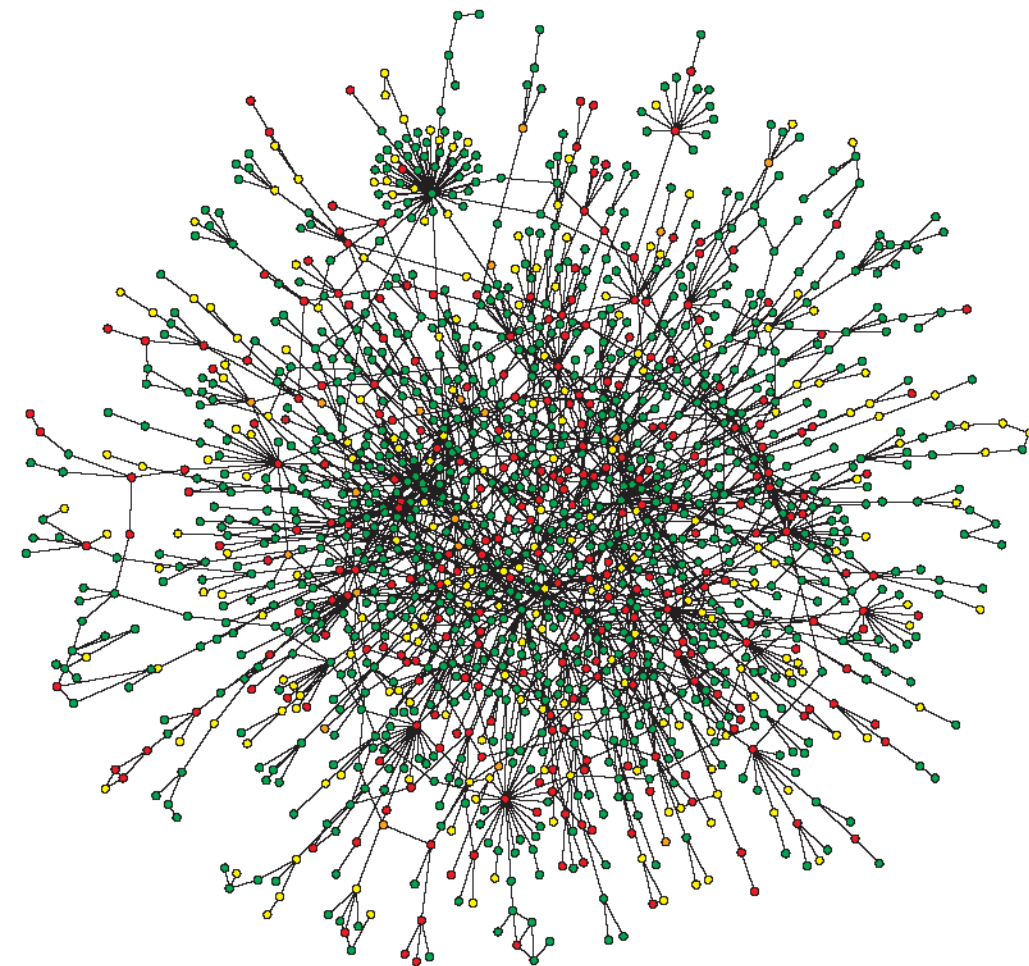
© 2016 Kelley et al.; Published by Cold Spring Harbor Laboratory Press

# ... предсказывает доступность хроматина В КЛЕТОЧНЫХ ЛИНИЯХ

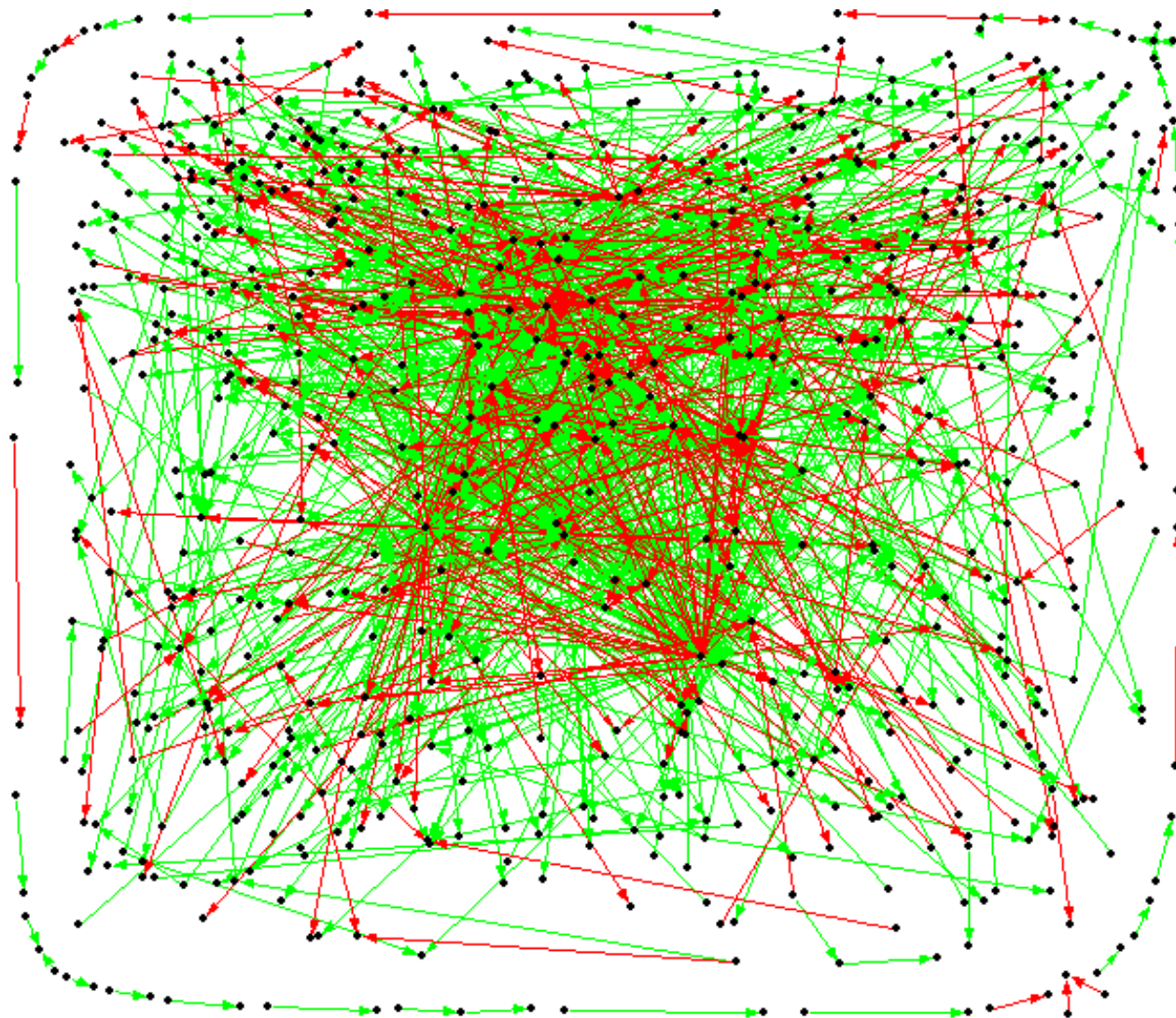




# Белок-белковые (структурные, сигнальные и др.) и белок-ДНКовые (регуляторные) взаимодействия в дрожжах



# Регуляция транскрипции у человека

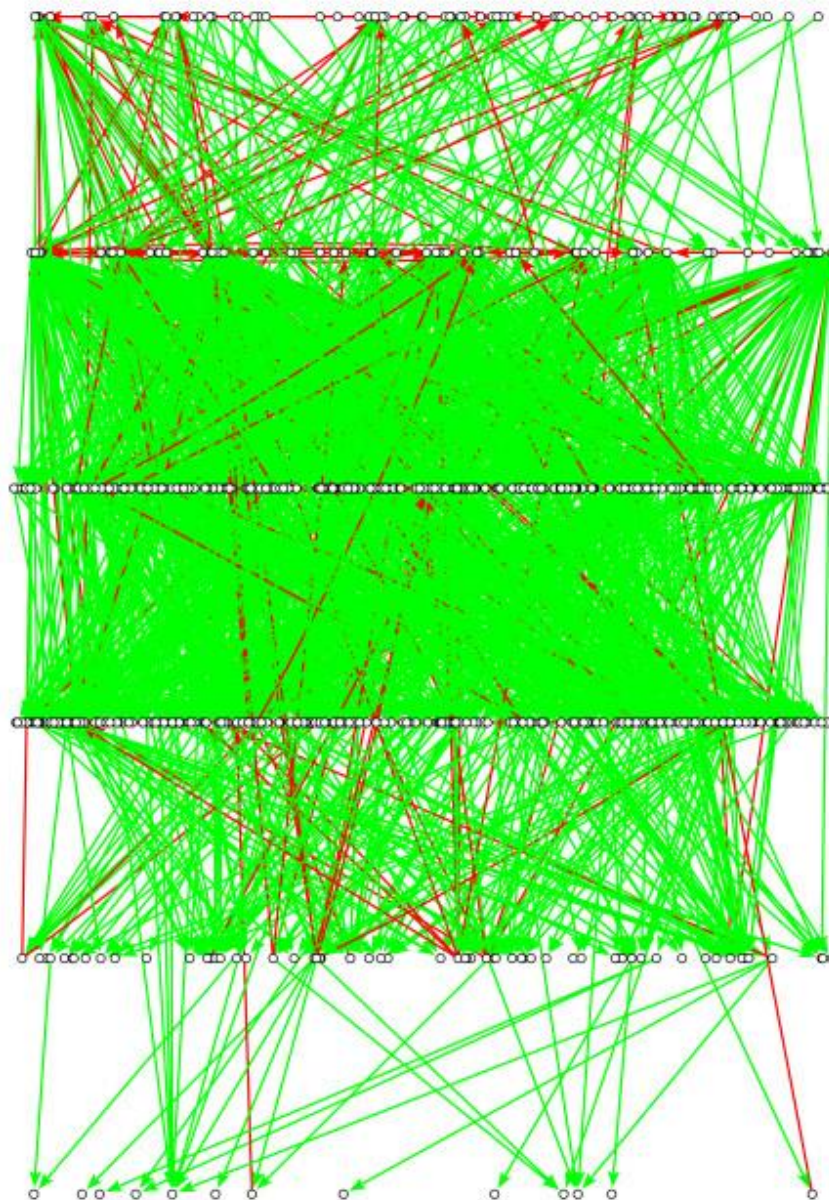


- 1449 взаимодействий между 689 генами
- Отношение «активаторы : репрессоры» = 3:1
- До 95 регулируемых генов, до 45 регуляторов.



**Иерархия:  
732 белков  
(71 рецепторов),  
1671 взаимодействий  
(фосфорилирование,  
дефосфорилирование,  
гидролиз etc)**

**208 анти-  
иерархических ребер**



# Динамика: активность транскрипционных взаимодействий в клеточных линиях

