**Join us at Mount Sinai!**
karr@mssm.edu
KarrLab.org

# Acknowledgements

# Outline

## Genotype to cellular phenotype

- What is a WC model?
- Why do we need WC models?
- Challenges & feasibility
- Foundational principles and state of the art
- Progress toward comprehensive models
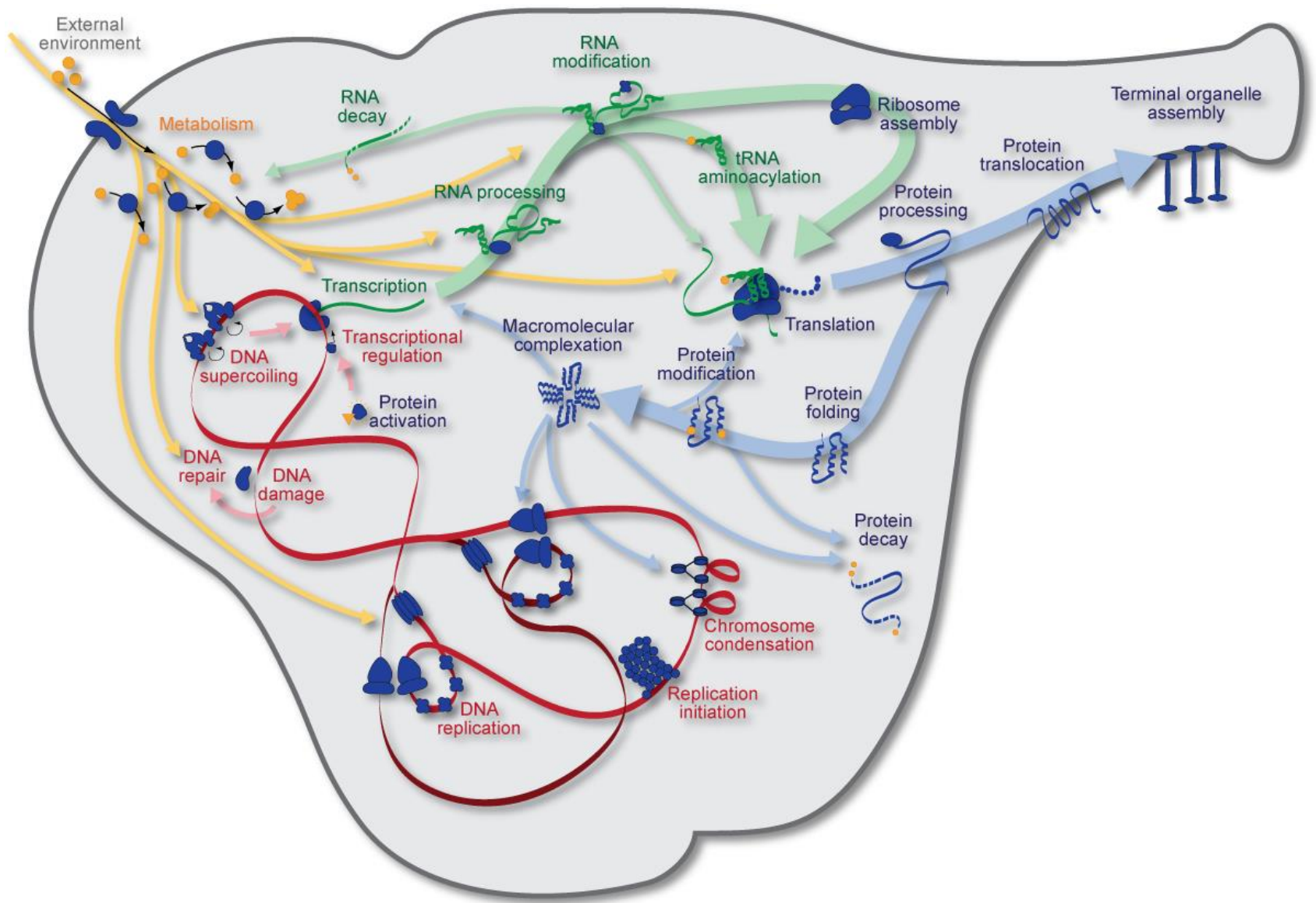
## Tips for modeling complex systems

# What is a WC model?

# Goals of WC modeling

**Species-specific**

**Whole cell**

**Whole genome**

**Whole cell cycle**

**Mechanistic**

**Dynamic**

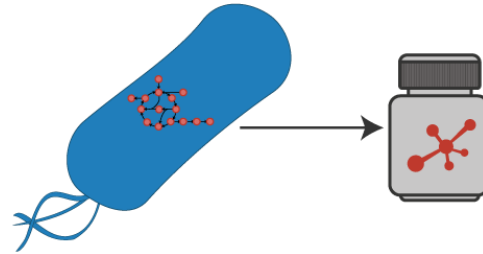**Stochastic**

AGTC

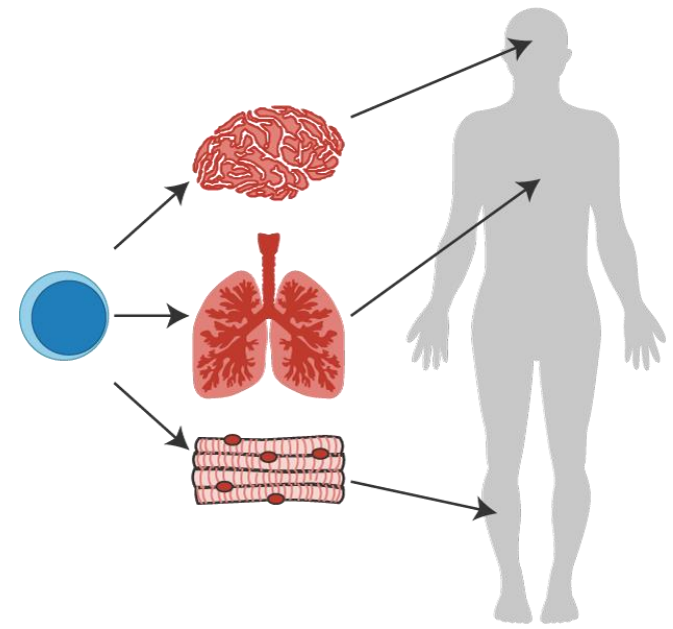Whole-Cell Modeling

Mount Sinai

# Motivation

# Synthetic biology requires WC models

Biosensors

Biofactories

Tissue engineering

# Example: drug biosynthesis

External environment

Metabolism

# Example: drug biosynthesis

# Challenges

# Challenge: explain diverse chemistry

Metabolism
*FBA*

Transcriptional
regulation
*Logical*

Signaling
*ODE, SSA*

# Challenge: explain multiple scales



Length

Growth

Replication

Transcription

Metabolism

Time

Whole-Cell Modeling

Mount Sinai

# Challenge: heterogeneous data



Transcription
*RNA-seq*



Protein expression
*Mass-spec, Western blot*



Single-cell variation
*Microscopy*

# Feasibility

Protein Abundance (0-12)

# Feasibility: Multi-algorithm simulation

**Genomic and biochemical data**

**Pathway submodels**

**Rule-based modeling**

**Multi-algorithmic simulation**

# Workflow

Data

Knowledge

Curate
databases →
& papers

Pathway/
genome
DB

Model

Experimental
conditions

Parallel
simulation

Visualization
& analysis

Whole-Cell
Modeling

Mount
Sinai

# 1. Focus on simple cells



| | ***E. coli*** | ***M. genitalium*** |
|---|---|---|
| **Genome** | 4700 kb | 580 kb |
| **Genes** | 4461 | 525 |
| **Size** | 2 μm × 0.5 μm | 0.2-0.3 μm |

Chemical reactions, catalysis, kinetics

Genetic sequence, genes, operons, promoters, transcriptional regulation

RNA processing, modification reactions

Signal sequences, localization, folding, chaperones, modification

Small molecules

DNA

RNA

Protein

Growth media composition

Metabolite structures, biomass composition

DNA-binding, motifs, sites, footprints

RNA expression, half-life, essentiality

Complex composition

Metabolism

Transcriptional
regulation

Signaling

Boolean
*Bolouri, 2000's*

FBA
*Palsson, 1990's*

ODE
*Shuler, 1970's*

PDE
Gillespie
*Luthey-Schulten, 2011*

Scope

Detail

# 3. Model each process

## Metabolism



## Species and reactions

ADP + Pi + 4 H+[p]

↓ ATPase

ATP + H2O + 3 H+[c]

## Catalysis

ATPase = 3 * AtpA
1 * AtpB
1 * AtpC
3 * AtpD
10 * AtpE
2 * AtpF
1 * AtpG
1 * AtpH

## Kinetics

$$v = \text{k}_{\text{cat}}[\text{ATPase}] \frac{[\text{ADP}]}{K_m + [\text{ADP}]}$$

States

Submodels

# 6. Verify model

## ☑ Matches training data

- ☑ Cell mass, volume
- ☑ Biomass composition
- ☑ RNA, protein expression, half-lives
- ☑ Superhelicity

## ☑ Matches published data

- ☑ Metabolite concentrations
- ☑ DNA-bound protein density
- ☑ Gene essentiality

## ☑ Matches theory

- ☑ Mass conservation
- ☑ Central dogma
- ☑ Cell theory
- ☑ Evolution

## ☑ No obvious errors

- ☑ Plot model predictions
- ☑ Manually inspect data
- ☑ Compare to known biology

Whole-Cell Modeling

Mount Sinai

State of the art

## Growth rate
Simulation: multiple selected

- Wild-type #1
- Wild-type #2
- Wild-type #3
- Wild-type #4
- Wild-type #5

Growth rate (cell s-1) vs Time (h)

## Cell shape
Simulation: Wild-type #1

126.3 aL
289 nm
5.9 fg

## Metabolism
Simulation: Wild-type #1

Hi
Lo

## Replication initiation – oriC DnaA Boxes
Simulation: Wild-type #1

oriC  R5  R4  R3  R2  R1

## DNA replication, protein occupancy, methylation, & damage
Simulation: Wild-type #1

Chromosome 1
- Mother
- Daughter
- Methylation
- Strand break
- ssDNA
- dsDNA
- Protein

Chromosome 2

## Mature protein monomer expression
Simulation: Wild-type #1

1
50001
100001
150001
200001
250001
300001
350001
400001
450001
500001
550001

Hi
Lo

00 : 02 : 25

*lacI*

# WC models help purpose drugs



M. genitalium
525 genes

Identify fragile genes
*Systems modeling*

Identify inhibitors
*Bioinformatics*

**Assess active & binding sites homology**
*Bioinformatics*

--TMEPLTEAYLFAAARTEHIS
--DMDIRTEAMLFAASRREHLV
DEVITDKAEVLMFYAARVQLVE
--DVEDHSVHLLFSANRWEQVP

**Assess affinity**
*Structural modeling*

Tmk
1 gene

# Limitations of the *Mycoplasma* model

- Represents one of the smallest bacteria

- Ignores several processes

- Mispredicts several phenotypes

- Methods were ad hoc

- Hard to understand, reuse, and expand

- Time-consuming to build

Whole-Cell Modeling

Mount Sinai

Toward more comprehensive and more accurate models

- Karyotypically normal

- Autonomous

- Well-characterized

# Bottlenecks



Data

Knowledge

Curate
databases
& papers

Pathway/
genome
DB

Model

Experimental
conditions

Parallel
simulation

Visualization
& analysis

Whole-Cell Modeling

Mount Sinai

# Bottlenecks

- **Data aggregation:** Hard to find relevant data
  - Data is incomplete, scattered, and insufficient annotated

- **Model design:** Hard to capture multiple scales and describe models modularly
  - Insufficient abstraction and metadata

- **Simulation**: Hard to simulate multiple scales
  - Simulators are only support individual formalisms and are slow

- **Verification:** Little formalism or standardization

- **Collaboration:** Difficult to describe the data, assumptions, and decisions that underlie modeling

$$v = k_{\text{cat}} \, [\text{enzyme}] \frac{[\text{substrate}]}{[\text{substrate}] + K_{\text{m}}}$$

Metabolite
concentrations

Enzyme
concentrations

Reaction kinetics

# Data needed for WC modeling

# Datanator: data integration & discovery



**Aggregate**

**Find**
- Species
- Environment

**Reduce**

$7.3 \cdot 10^{-4}$ mM

**Review**

kineticdatanator.com

[Phosphate]

6.3 E -2 mM

# Datanator: data aggregation

**Metabolites**
- ChEBI
- ECMDB, YMDB
- PubChem

**DNA**
- GenBank

**RNA**
- Array Express
- MODOMICS
- RNALocate
- RNA MOD

**Protein**
- COMPARTMENTS
- CORUM
- Human Protein Ref. DB
- Pax-DB
- PDB
- PSORTdb
- RESID
- UniProt

**Interactions**
- BioCyc
- DBTBS
- DrugBank
- JASPAR
- KEGG
- SuperTarget

**Taxonomy**
- NCBI

**Pathways**
- KEGG
- Pathway Commons
- Reactome
- WikiPathways

**Rates**
- BRENDA
- SABIO-RK

Whole-Cell Modeling

Mount Sinai

# Datanator: actionable metadata

**Measured entity/property**

**Measured value, uncertainty, units**

**Genotype**
– Taxon
– Genetic variant
– Cell, tissue type

**Environment**
– Temperature
– pH
– Growth media

**Data generation process**
– Experimental design
– Measurement method

**Data analysis process**
– Software
– Version

**Metadata**
– Authors
– Curator
– Date
– Citation

Whole-Cell Modeling

Mount Sinai

## Chemical similarity

– Tanimoto index

– Sequence similarity

## Genetic similarity

– Whole-genome similarity

– Taxonomic distance

## Environmental similarity

– Temperature

– pH

# WC-Lang: scalable model descriptions

- Concretely describe composite multi-algorithmic models

- Concrete descriptions of every model element

- Capture data and assumptions underlying models

- Explicit descriptions of mixed granularity / lumping

- Structured description of initial conditions

- User interfaces suited to large models

Whole-Cell Modeling

Mount Sinai

RNA(i, 0) + NTP(i, 1) $\rightarrow$ RNA(i, 1) + PPi

RNA(i, 1) + NTP(i, 2) $\rightarrow$ RNA(i, 2) + PPi

RNA(i, 2) + NTP(i, 3) $\rightarrow$ RNA(i, 3) + PPi

RNA(i, 3) + NTP(i, 4) $\rightarrow$ RNA(i, 4) + PPi

...

Whole-Cell Modeling

Mount Sinai

$$RNA(i, l) + NTP(i, l+1) \rightarrow RNA(i, l+1) + PPi$$

$$RNA(I, l) + H_2O \rightarrow RNA(i, l-1) + NMP(i, l)$$

$$Protein(i, l) + AA(i, l+1) \rightarrow Protein(i, l+1) + H_2O$$

$$Protein(i, l) + H_2O \rightarrow Protein(i, l-1) + AA(i, l)$$

# WC-Lang: scalable model descriptions

|          | Initiation | | Elongation | | Termination | |
|----------|------------|-----|------------|--------|-------------|-----|
| SBML     | 1 per RNA  | 335 | 1 per base | ~500k  | 1 per RNA   | 335 |
| Rules    | 1          | 1   | 1          | 1      | 1           | 1   |

# H1-hESC model



Recon 2.2 → H1 model → Composable model

- H1 transcriptomics data (ENCODE)
- Cell composition
- Media composition

- Kinetic data (SABIO-RK)
- Protein abundance (Phanstiel et al., 2011; PaxDB)

# Summary

# Availability

- **Code:** code.karrlab.org (GitHub, PyPI)

- **Data:** data.karrlab.org (Quilt)

- **Images:** DockerHub

- **Primer and docs:** docs.karrlab.org

- **Tutorials:** sandbox.karrlab.org

Whole-Cell Modeling

Mount Sinai

# Summary

Bioengineering and medicine needs WC models

WC modeling is becoming feasible

New technologies will enable WC modeling

Pilot models will show the feasibility of bacteria and human models

Whole-Cell Modeling

Mount Sinai

Tips & tricks

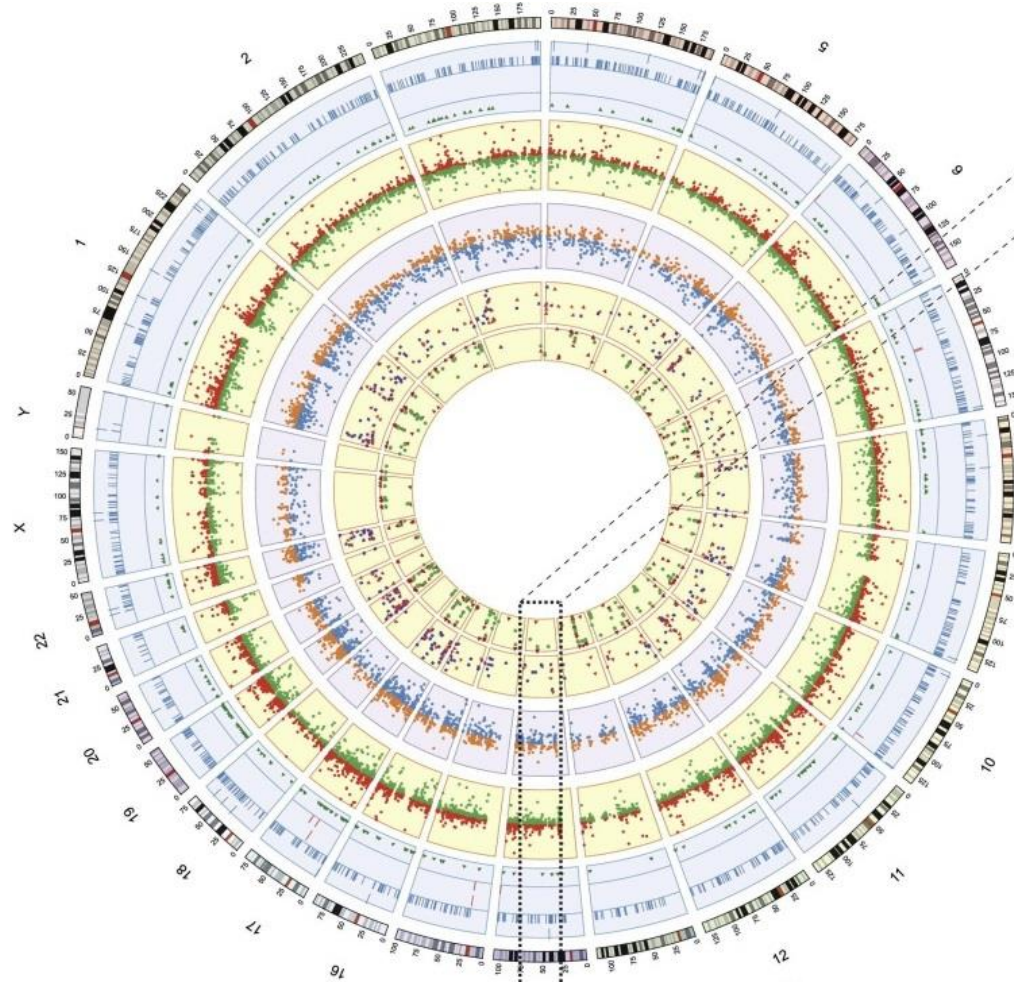# Integration enables great scope and depth

- Data aggregation

- Model composition

- Multi-algorithmic co-simulation

- Modular methods and software

- Interdisciplinary collaboration

# Frameworks enable scalable integration

# Common languages enable frameworks

A

# Collaboration enables solutions

# Modularity enables collaboration



Data

Knowledge

Curate
databases
& papers

Pathway/
genome
DB

Model

Experimental
conditions

Parallel
simulation

Visualization
& analysis

Whole-Cell
Modeling

Mount
Sinai

# Sharing promotes collaboration

- **Quilt:** data

- **GitHub:** code

- **PyPI:** packaged code

- **Docker:** computing environments

- **Google Docs, Overleaf:** written documents

- **Google Drive:** other files

- **GitHub issues:** tasks

Whole-Cell Modeling

Mount Sinai

# Common practices ease collaboration

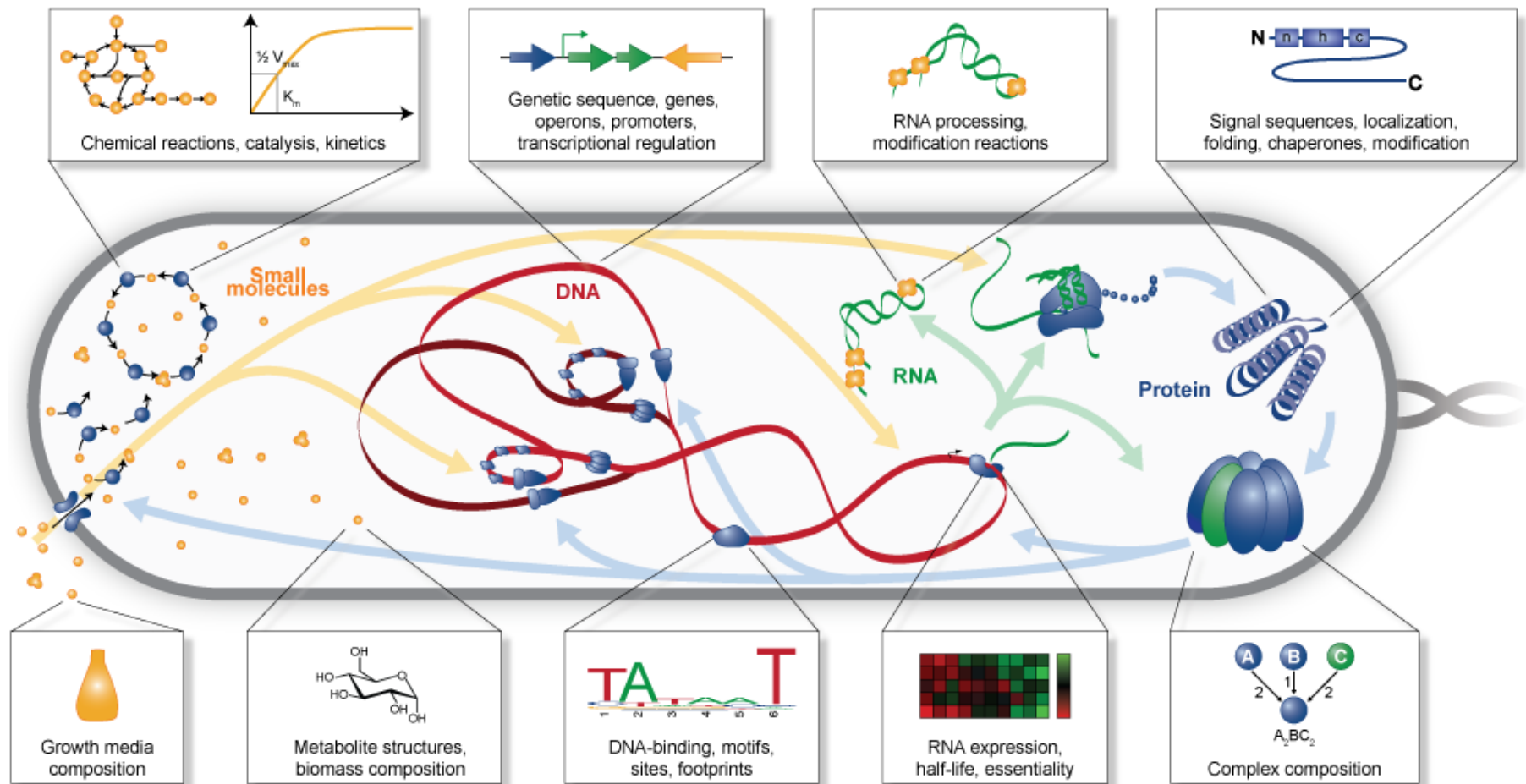- Interfaces between modules

- Coarse-graining

- Package organization

- Coding, documentation styles

- Software libraries

Whole-Cell Modeling

Mount Sinai

# QC inspires trust among collaborators

Chemical reactions, catalysis, kinetics

Genetic sequence, genes, operons, promoters, transcriptional regulation

RNA processing, modification reactions

Signal sequences, localization, folding, chaperones, modification

Small molecules

DNA

RNA

Protein

Growth media composition

Metabolite structures, biomass composition

DNA-binding, motifs, sites, footprints

RNA expression, half-life, essentiality

Complex composition

# Summary

# Integration is enabling WC modeling