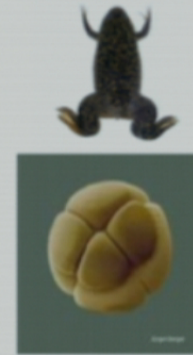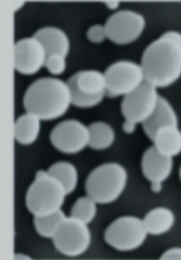Title: Inferring the spatiotemporal DNA replication program from noisy data

Date: Dec 05, 2013  02:40 PM

URL: http://pirsa.org/13120017

Abstract: <span>In eukaryotic organisms, DNA replication is initiated at â€œorigins,â€• launching â€œforksâ€• that spread bidirectionally to replicate the genome.  The distribution and firing rate of these origins and the fork progression velocity form the â€œreplication program.â€•  With Antoine Baker, I generalize a stochastic model of DNA replication to allow for space and time variations in origin-initiation rates, characterized by a function I(x,t).  We then address the inverse problem of inferring I(x,t) from experimental data concerning replication in cell populations.  Previous work based on curve fitting depended on arbitrarily chosen functional forms for I(x,t), with free parameters that were constrained by the data.  We introduce a model-free, non-parametric method of inference that is based on Gaussian process regression, a well-known inference technique from the machine-learning community.  The method replaces specific assumptions about the functional form of the initiation rate with more general prior expectations about the smoothness of variation of this rate, along the genome and in time.  Using this inference method, we can recover simulated replication schemes with data that are typical of current experiments without having to know or guess the functional form for the initiation rate I(x,t).  I will argue that Gaussian process regression has many other potential applications to physics.<br></span>

# Overview of DNA replication: Some numbers



|  | E. Coli (prok.) | S. Cerevisiae (euk.) | Xenopus (euk.) |
|---|---|---|---|
| DNA length | $4 \times 10^6$ bp | $12 \times 10^6$ bp | $3 \times 10^9$ bp (sperm) |
| Fork velocity | 1000 bp/s | 30 bp/s | 20 bp/s |
| S Phase (replication) | 25 min | 60 min | 20 min |
| Number of origins | 1 | 300 | $10^5$ |
| Sequence dependence? | Yes (oriC - 245 bp) | Yes (ACS) | No (early embryos) |

Images: *E. coli* (D. Kincaid), *S. cerevisiae* (microbiologyonline.org.uk), *Xenopus* (NIGMS and Uni. Tübingen)

# Overview of DNA replication: Timing-curve measurements

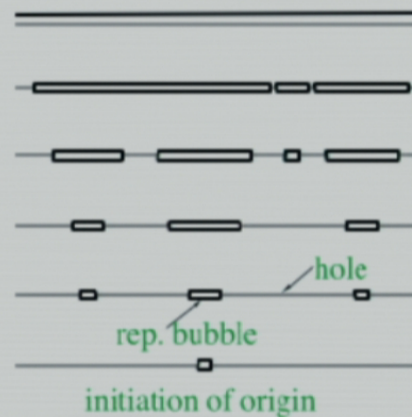McCune et al., Genetics 2008



Results:    Local replication fraction at a few times
            averaged over cell population
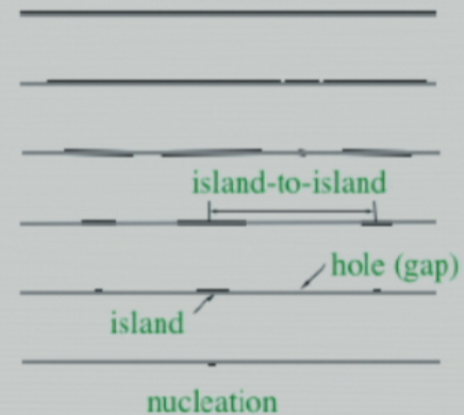Old way:    Microarrays
New way:    Sequencing

# Replication theory: Kinetic model

## DNA replication (Biology)



- hole
- rep. bubble
- initiation of origin

- Initiation of origin
- Replicated domains
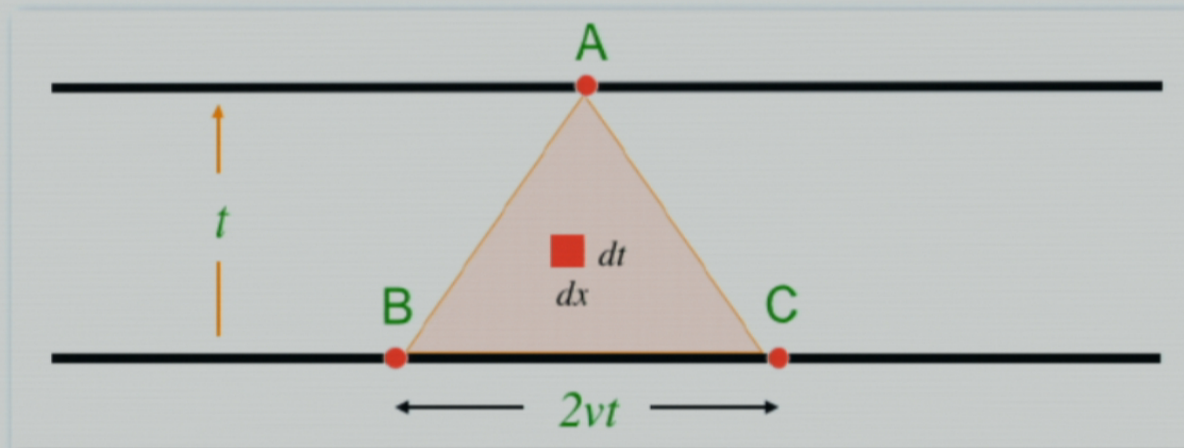- Unreplicated domains
- Fork velocity

## Crystal growth (Physics)

Time



- island-to-island
- hole (gap)
- island
- nucleation

- Nucleation of crystal
- Solid
- Liquid
- Growth velocity

3D: Kolmogorov; Johnson & Mehl; Avrami (1930s)
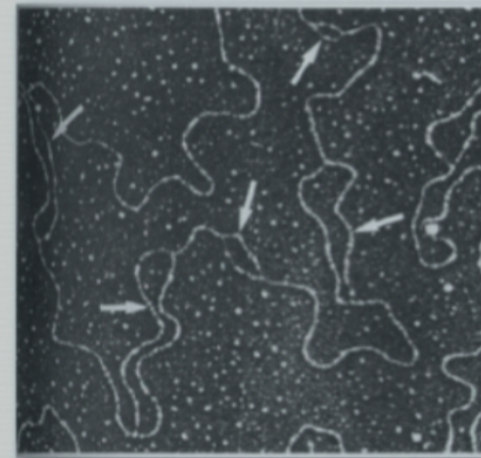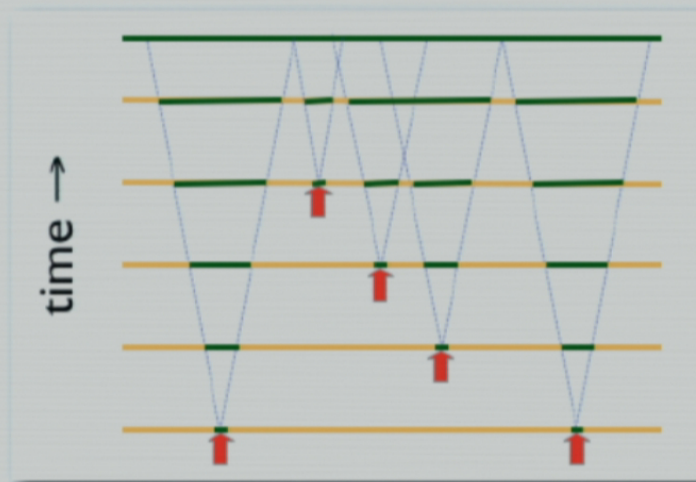ID: Sekimoto; Ben-Naim, Krapivsky (1980s/90s)

# Replication theory: Kinetic model



$s(x,t)$ = probability of **not** being replicated at *x* at time *t*  (Poisson)

$$s(x,t) = \prod_{\triangle}[1 - I(x',t')\,dx'\,dt'] = e^{-\iint_{\triangle} I(x',t')\,dx'\,dt'}$$
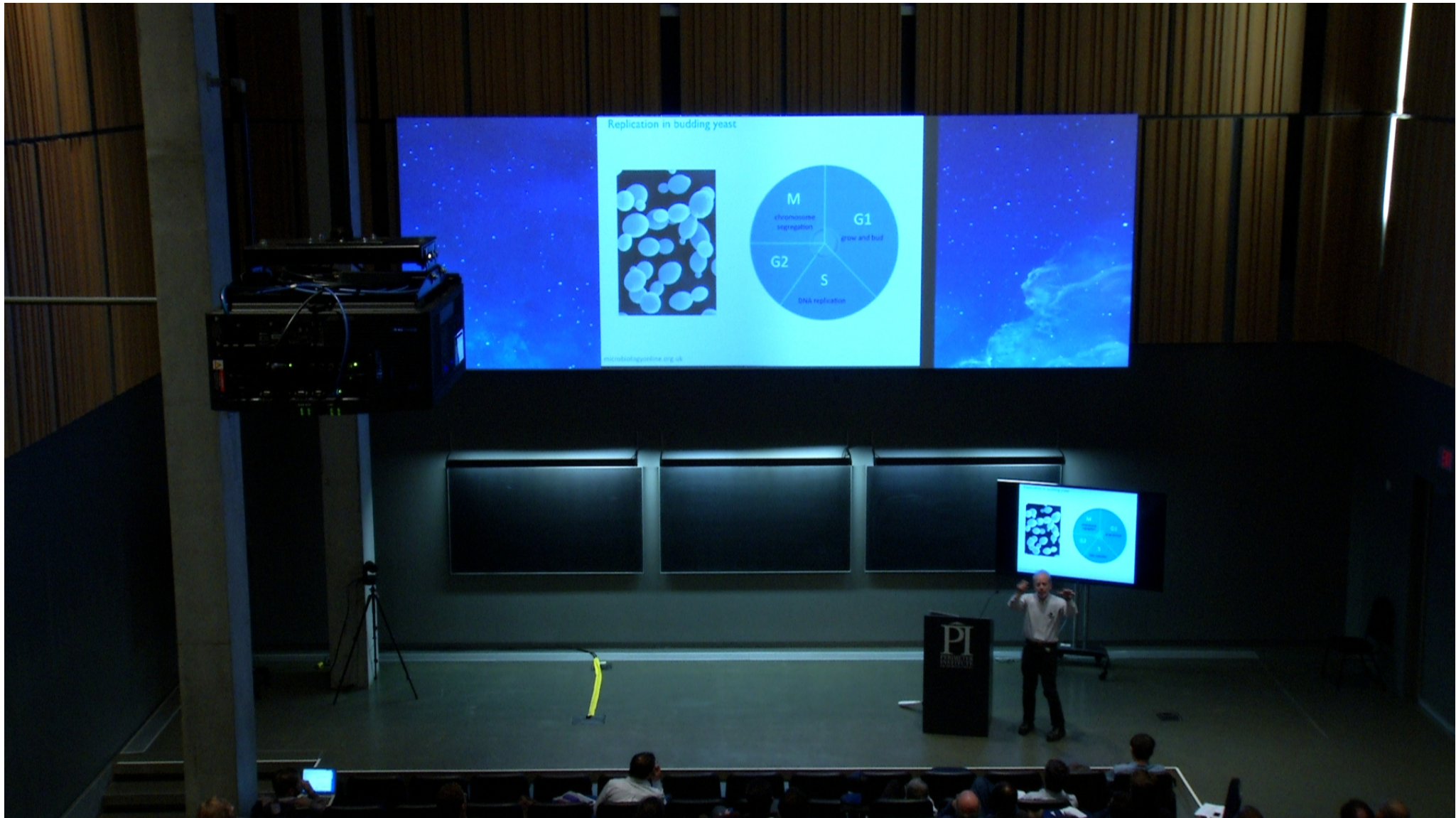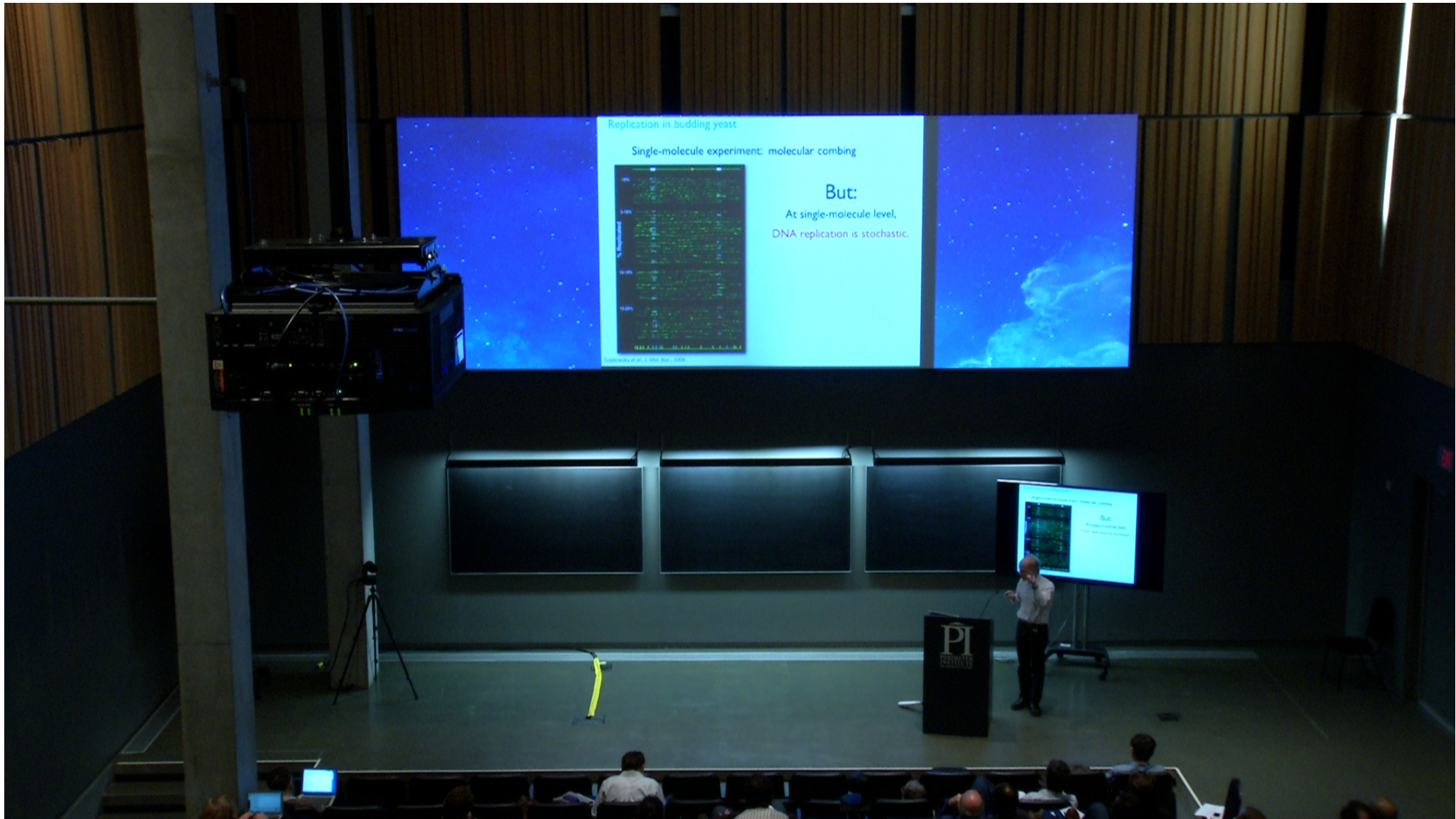
Kolmogorov, 1937
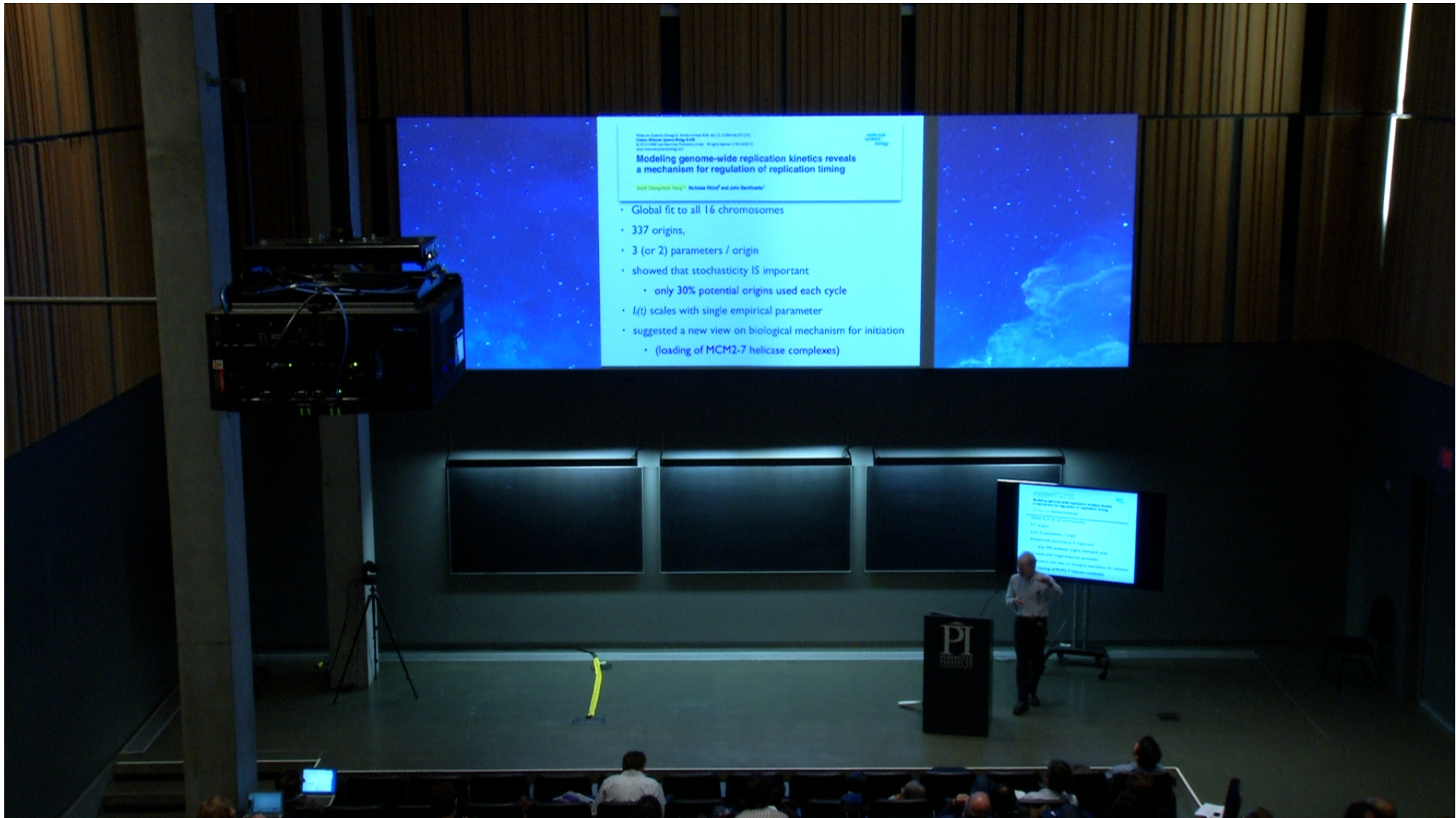
Replication theory: Kinetic model

Initiation rate $I(x,t)$ & growth rate $v$ are unknown in DNA replication

time →

Extract $I(x,t)$ and $v$
(spatiotemporal program of DNA replication)
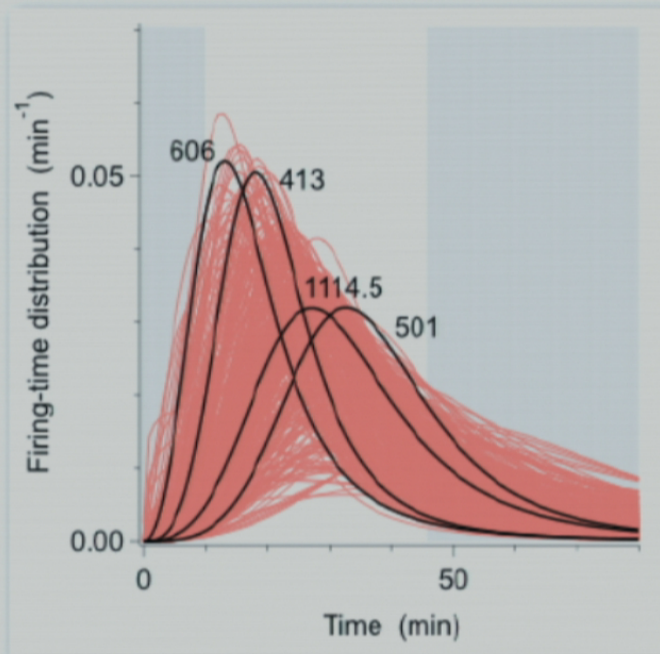
molecular
systems
biology

## Modeling genome-wide replication kinetics reveals a mechanism for regulation of replication timing

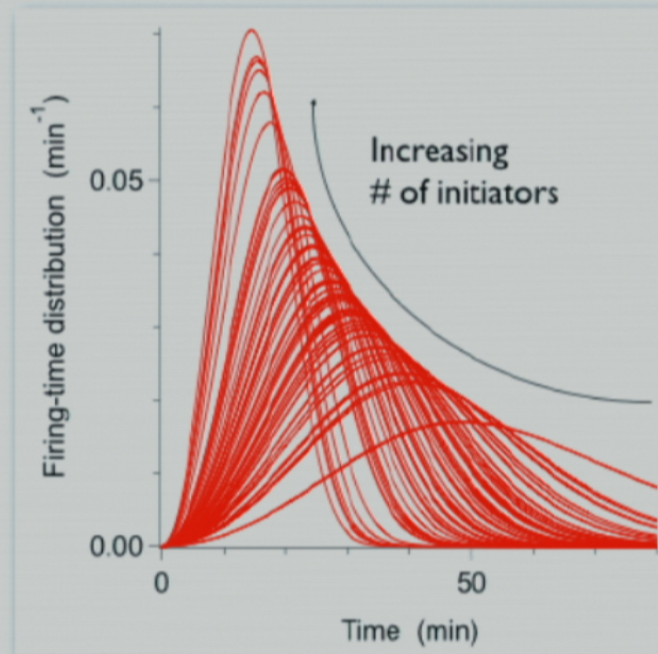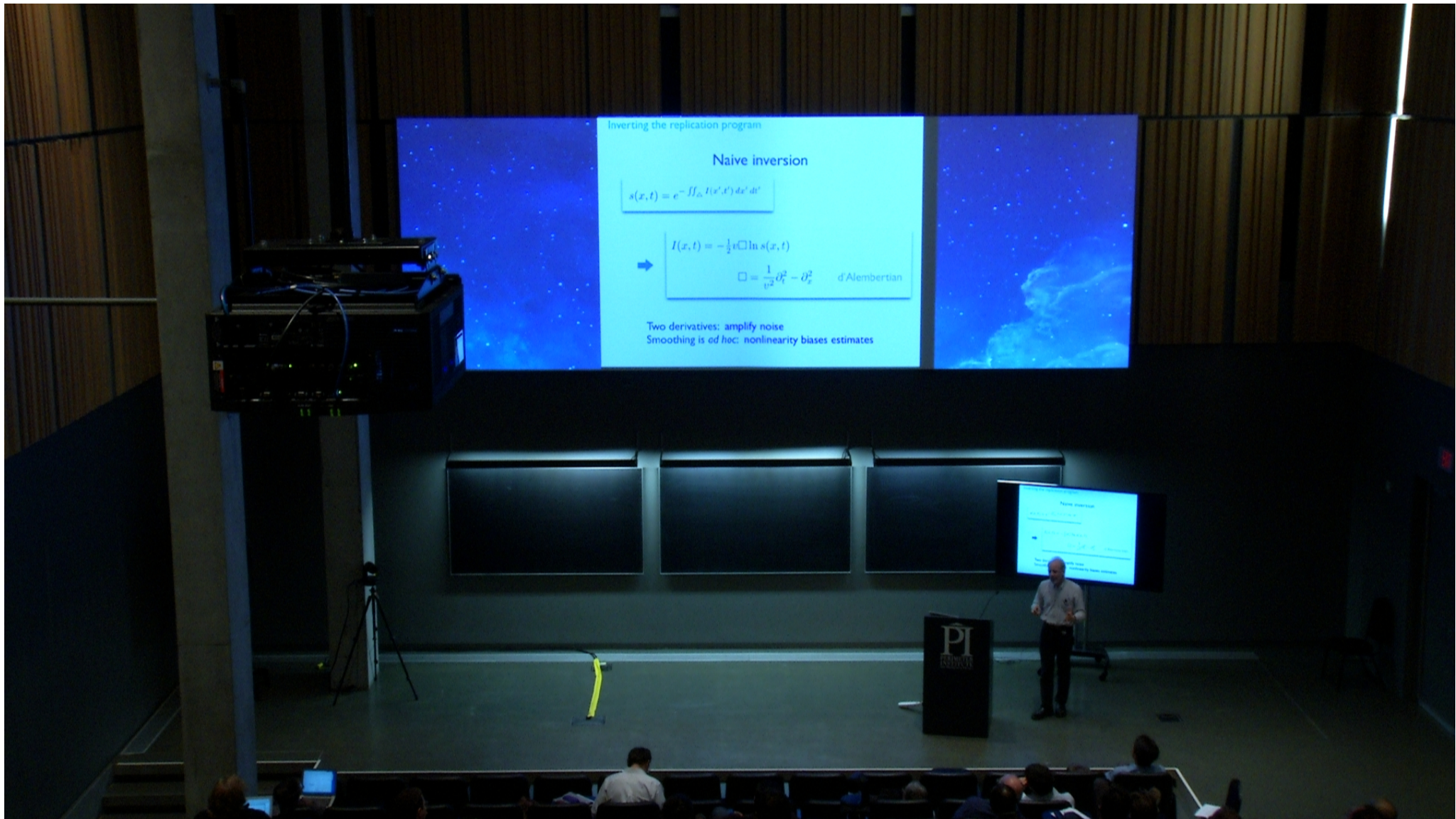Scott Cheng-Hsin Yang[1,a], Nicholas Rhind[2] and John Bechhoefer[1]

- Global fit to all 16 chromosomes

- 337 origins,

- 3 (or 2) parameters / origin

- showed that stochasticity IS important

  - only 30% potential origins used each cycle

- $I_i(t)$ scales with single empirical parameter

- suggested a new view on biological mechanism for initiation

  - (loading of MCM2-7 helicase complexes)

Modeling genome-wide replication kinetics reveals
a mechanism for regulation of replication timing

Scott Cheng-Hsin Yang[1,*], Nicholas Rhind[2] and John Bechhoefer[1]

*Tour de force*,  BUT:

- needed a talented physics student to pull it off

- initial guesses for nonlinear curve fit

- unknown number of parameters

- budding yeast is one of the best-studied cases

- genome only 12 MB long  (human = 3000 MB)