

# Thoughts on Publication Rollout Structure for the ENCODE Project

M Gerstein

(Publication analysis done by  
D Wang, KK Yan, J Rozowsky, E Pan)

The Human Genome Project



ENCODE Pilot



ENCODE Production



Comparative ENCODE



Epigenome Roadmap

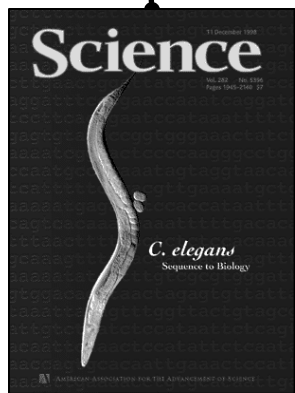


2000

2005

2010

2015



Worm Genome



modENCODE



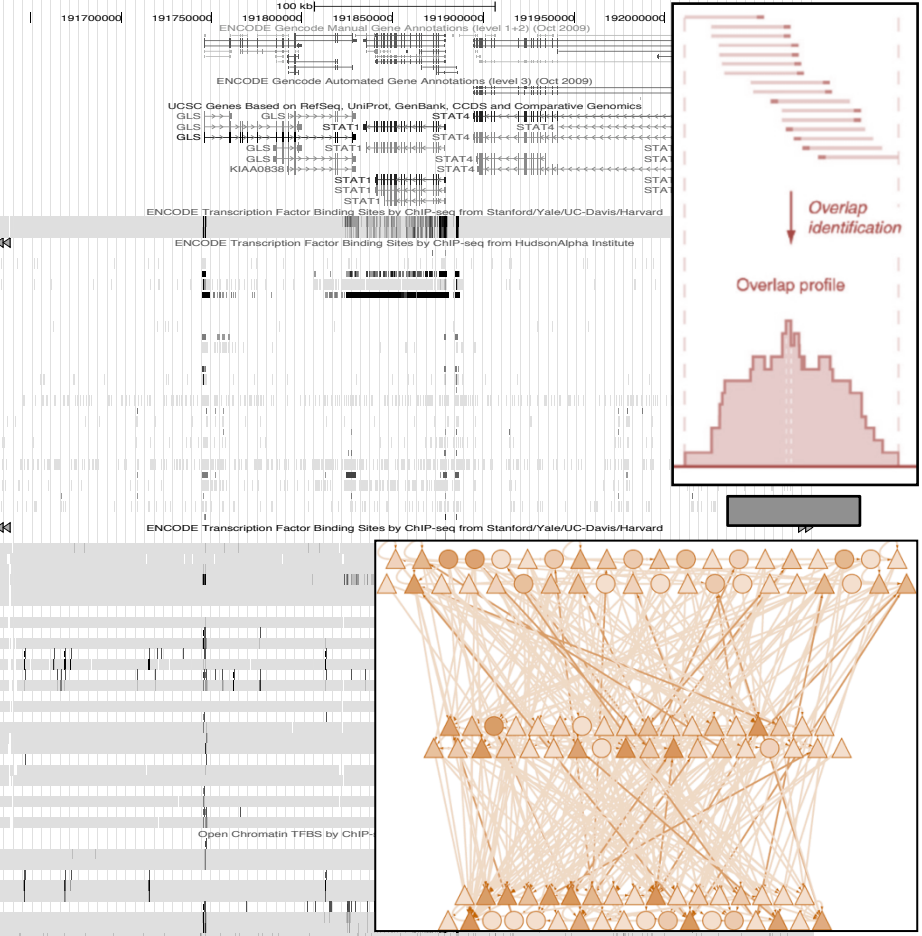
1000 Genomes Pilot



1000 Genomes Phase 3



GTEx

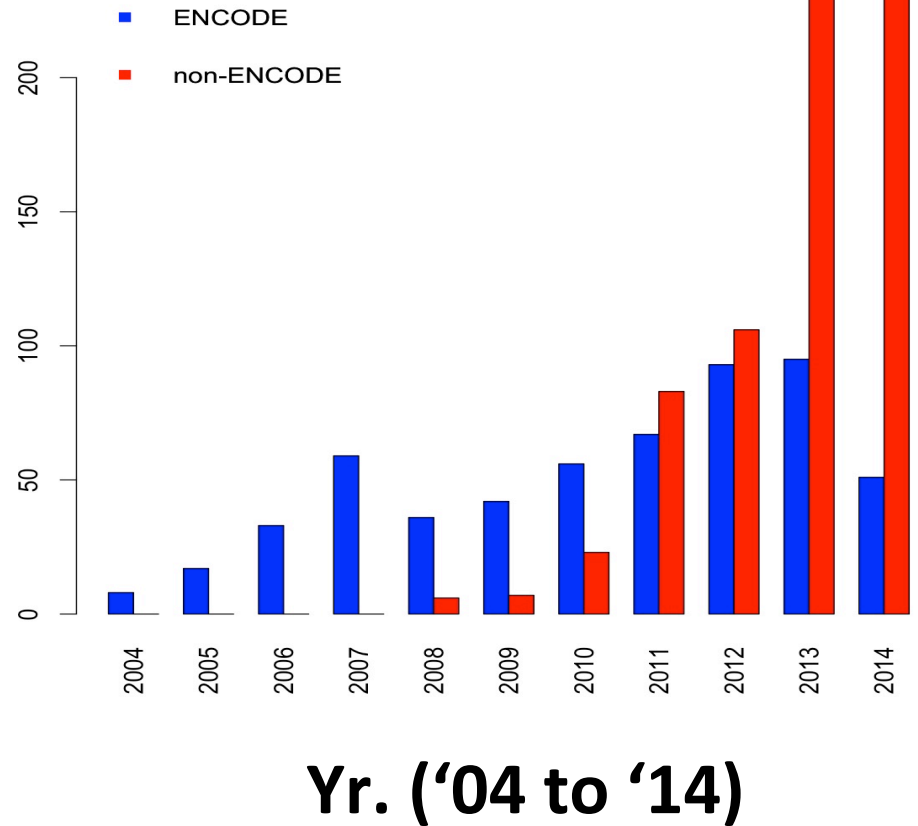
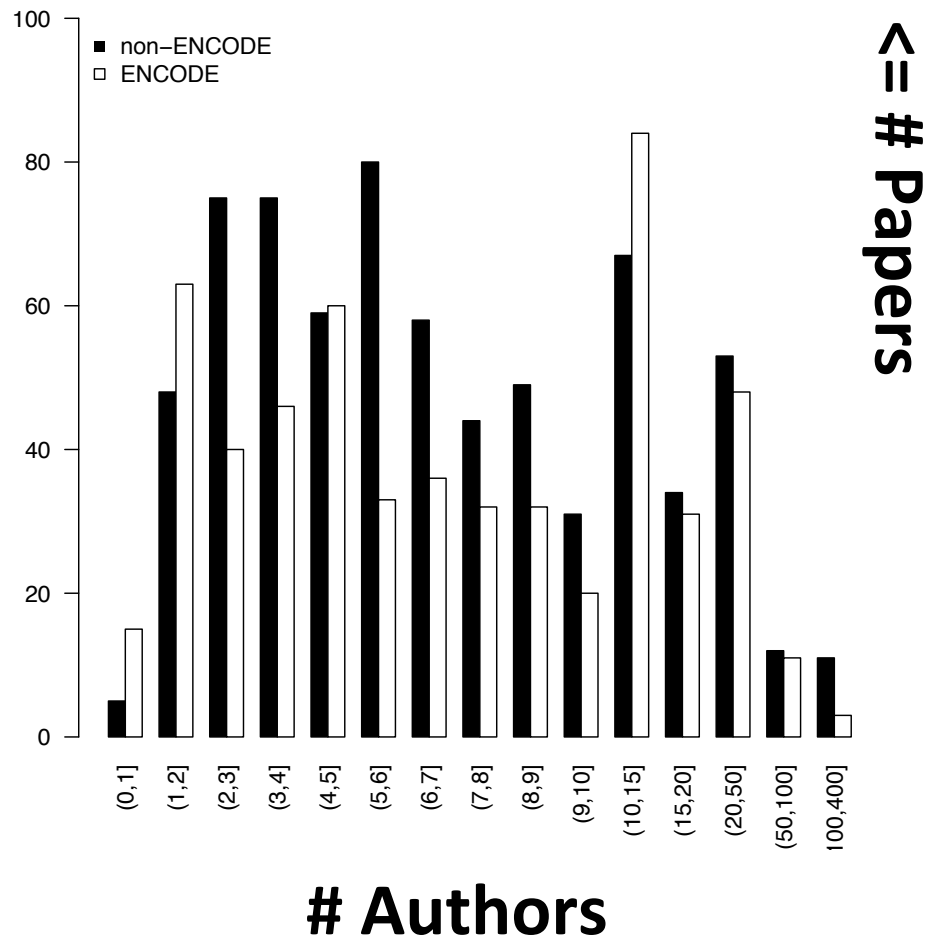


Overall ENCODE presentation as a **Structured Hierarchy** of huge amount of Genomic Data

- **Raw data** (reads) at the bottom
- Progressive Processed **Summaries**
  - Signals (e.g. how much DNA is bound by TFs)
  - Site locations
  - Reg. networks, chromatin states & stat. models
- **Code & VMs + Result Stats**
- **Many linked publications** are near top, documenting everything & forming metadata
- **Abstract of Consortium** paper sits at pinnacle



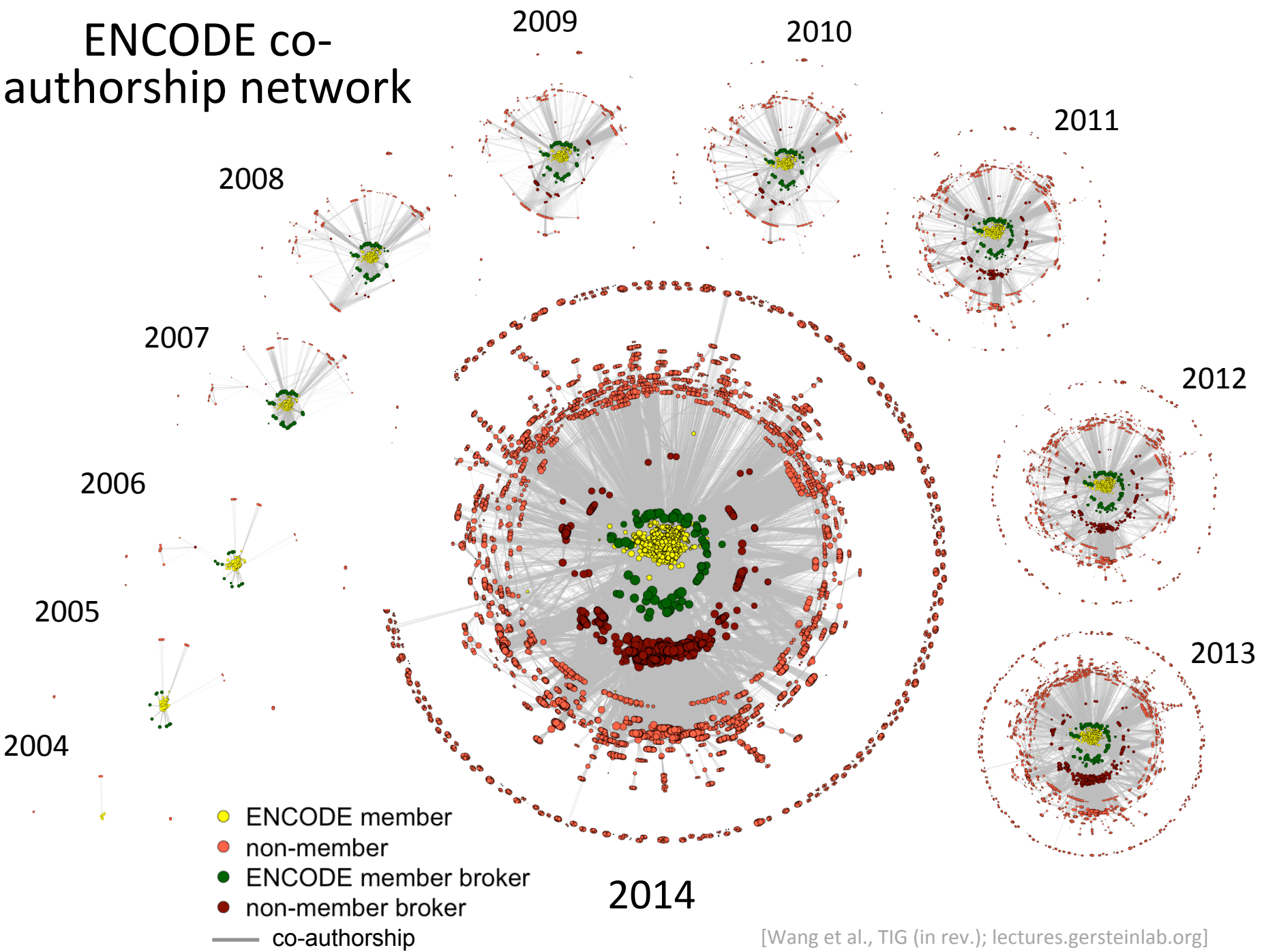
Reader/User goes **Top-Down**  
 Creator/Author builds **Bottom-up**



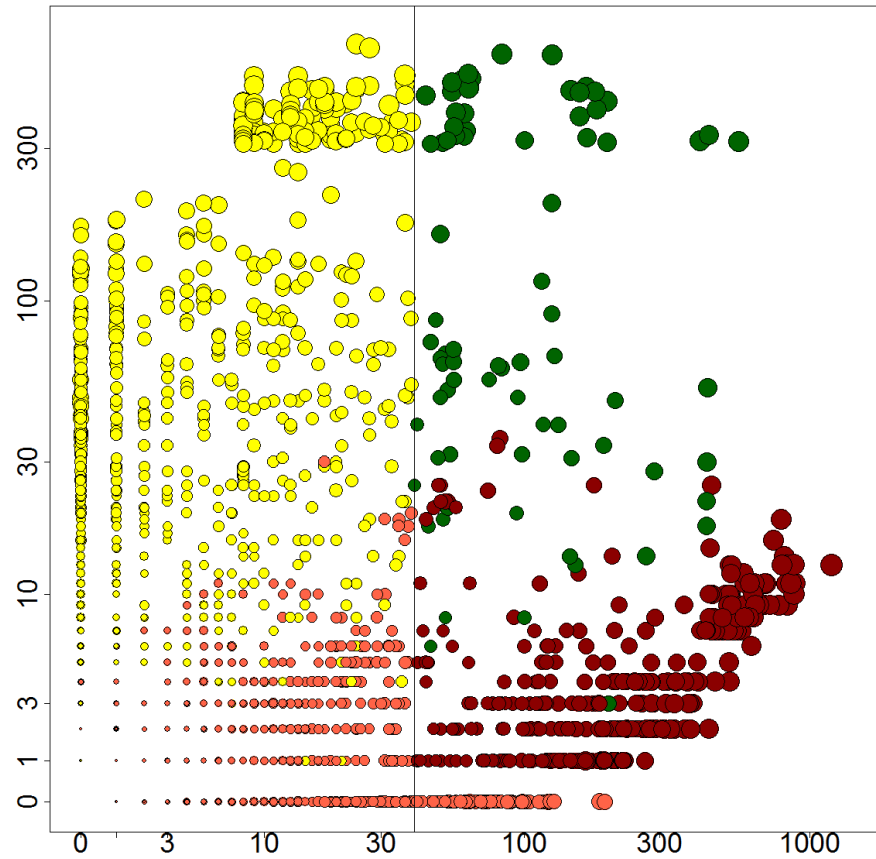
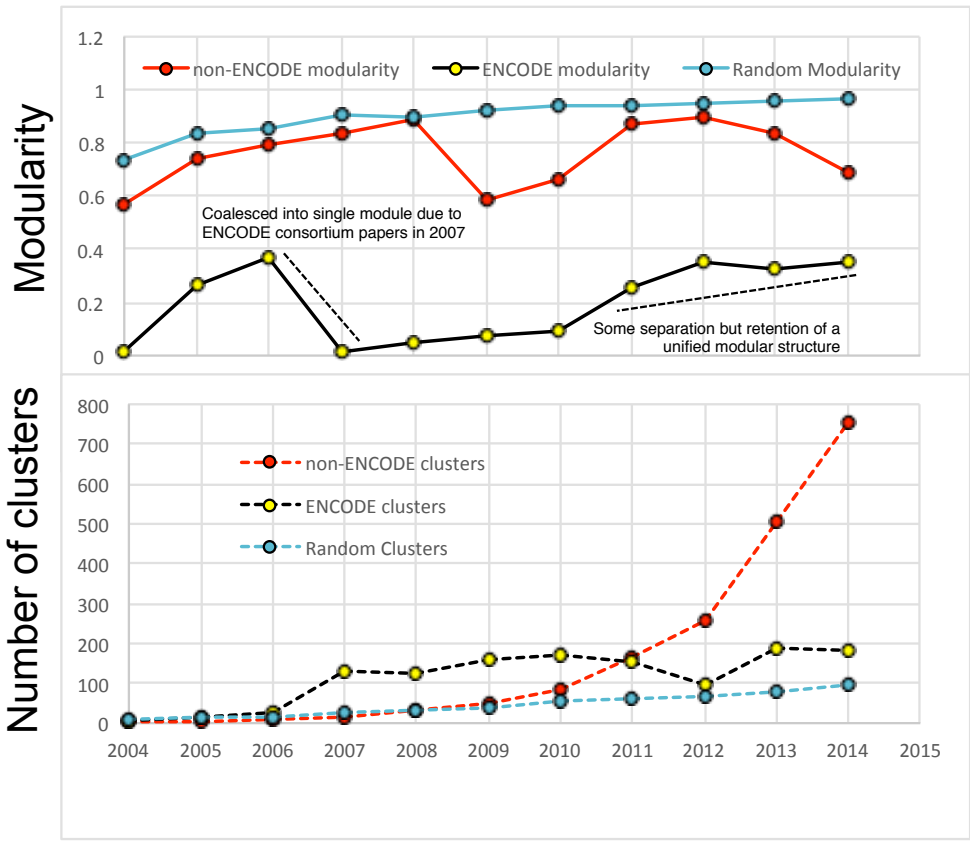
Papers authored by ENCODE consortium members vs. those that use ENCODE data but were not funded by ENCODE



# ENCODE co-authorship network



# Network statistics highlight change in modularity with consortium rollouts (L) & importance of broker role (R)



- ENCODE member
- non-member
- ENCODE member broker
- non-member broker