

# Understanding Random SAT

## Beyond the Clauses-to-Variables Ratio

Eugene Nudelman  
*Stanford University*

*joint work with...*

Kevin Leyton-Brown  
Holger Hoos

*University of British Columbia*

Alex Devkar  
Yoav Shoham  
*Stanford University*

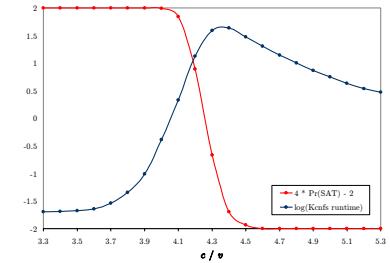
# Introduction

- SAT is one of the **most studied** problems in CS
- Lots known about its **worst-case** complexity
  - But often, particular instances of  $\mathcal{NP}$ -hard problems like SAT are **easy in practice**
- “Drosophila” for **average-case** and **empirical** (typical-case) complexity studies
- (Uniformly) random SAT provides a way to bridge analytical and empirical work



# Previously...

- **Easy-hard-less hard** transitions discovered in the behaviour of DPLL-type solvers [Selman, Mitchell, Levesque]
  - Strongly correlated with phase transition in solvability
  - Spawned a new enthusiasm for using empirical methods to study algorithm performance
- Follow up included study of:
  - Islands of tractability [Kolaitis et. al.]
  - SLS search space topologies [Frank et.al., Hoos]
  - Backbones [Monasson et.al., Walsh and Slaney]
  - Backdoors [Williams et. al.]
  - Random restarts [Gomes et. al.]
  - Restart policies [Horvitz et.al, Ruan et.al.]
  - ...

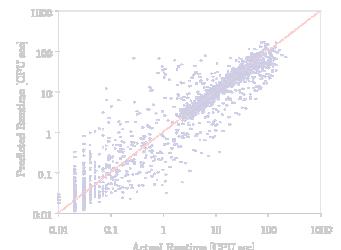
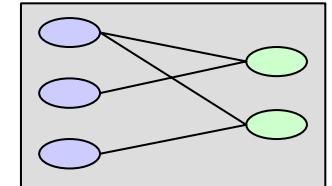


# Empirical Hardness Models

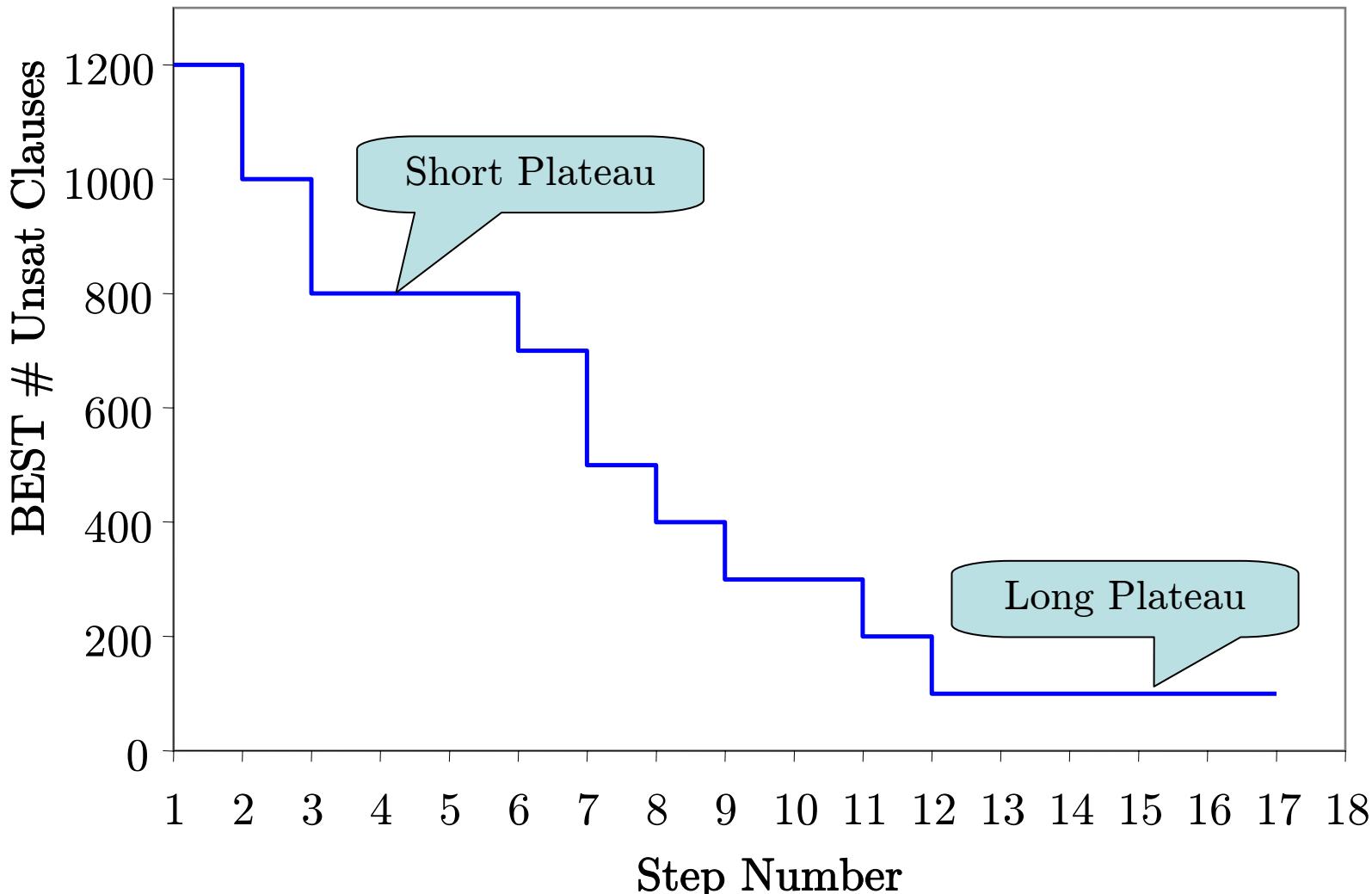
- We proposed building **regression models** as a disciplined way of predicting and studying algorithms' behaviour  
[Leyton-Brown, Nudelman, Shoham, CP-02]
- **Applications** of this machine learning approach:
  - 1) Predict running time
    - Useful to know **how long** an algorithm will run
  - 2) Gain theoretical understanding
    - Which variables are **important** to the hardness model?
  - 3) Build algorithm portfolios
    - Can select the right algorithm on a **per-instance** basis
  - 4) Tune distributions for hardness
    - Can generate **harder** benchmarks by rejecting easy instances

# Outline

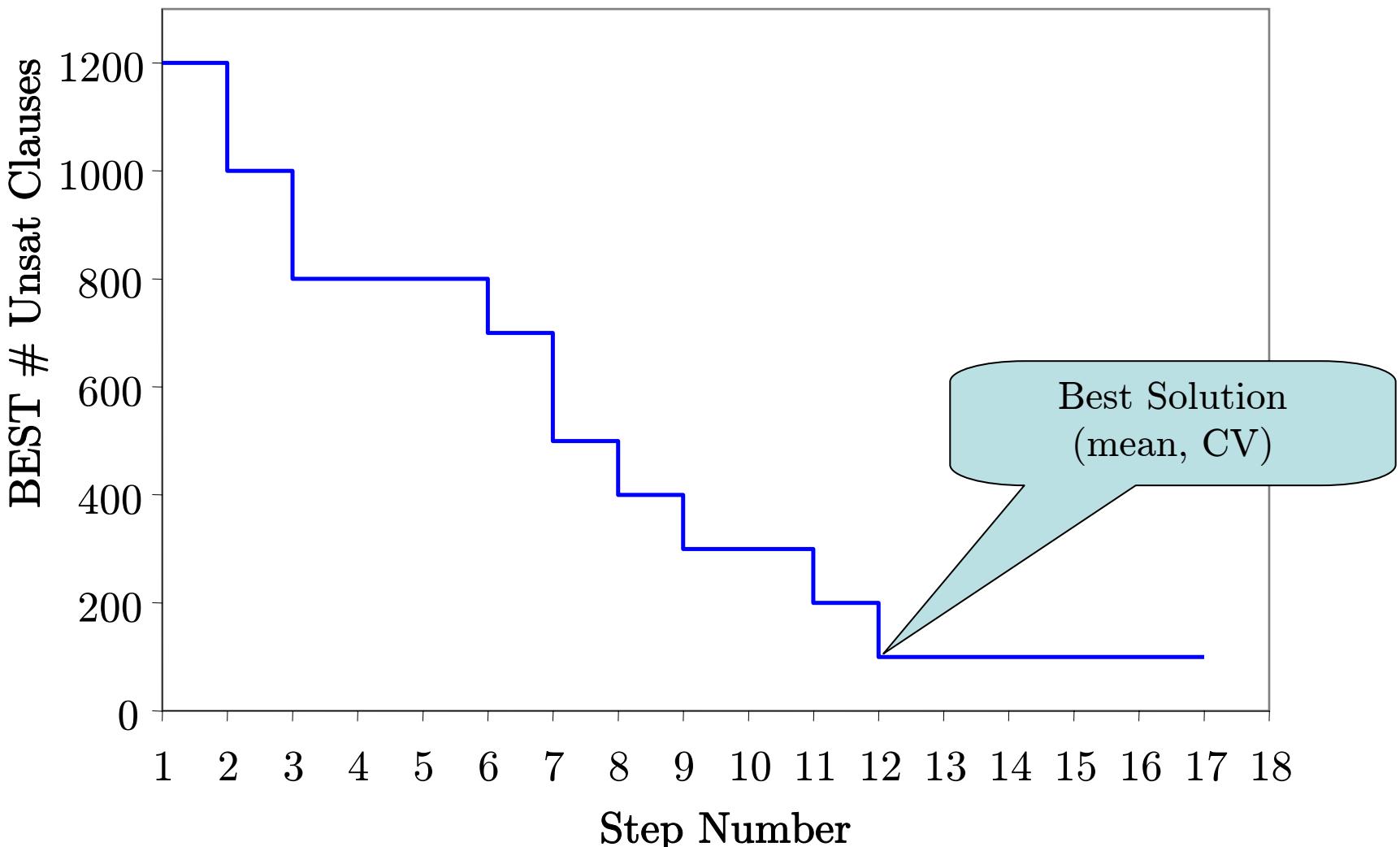
- Features
- Experimental Results
  - Variable Size Data
  - Fixed Size Data



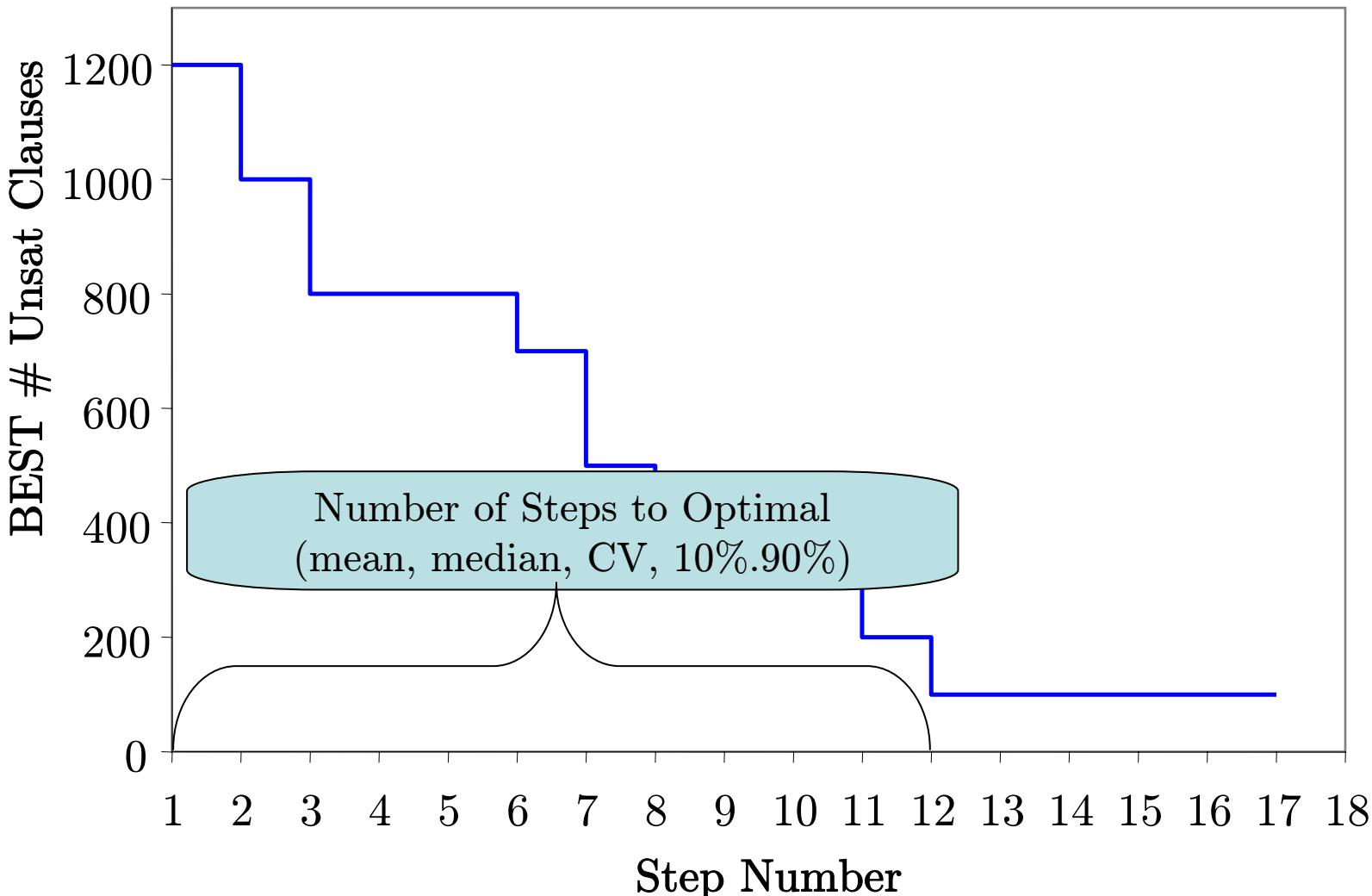
# Features: Local Search Probing



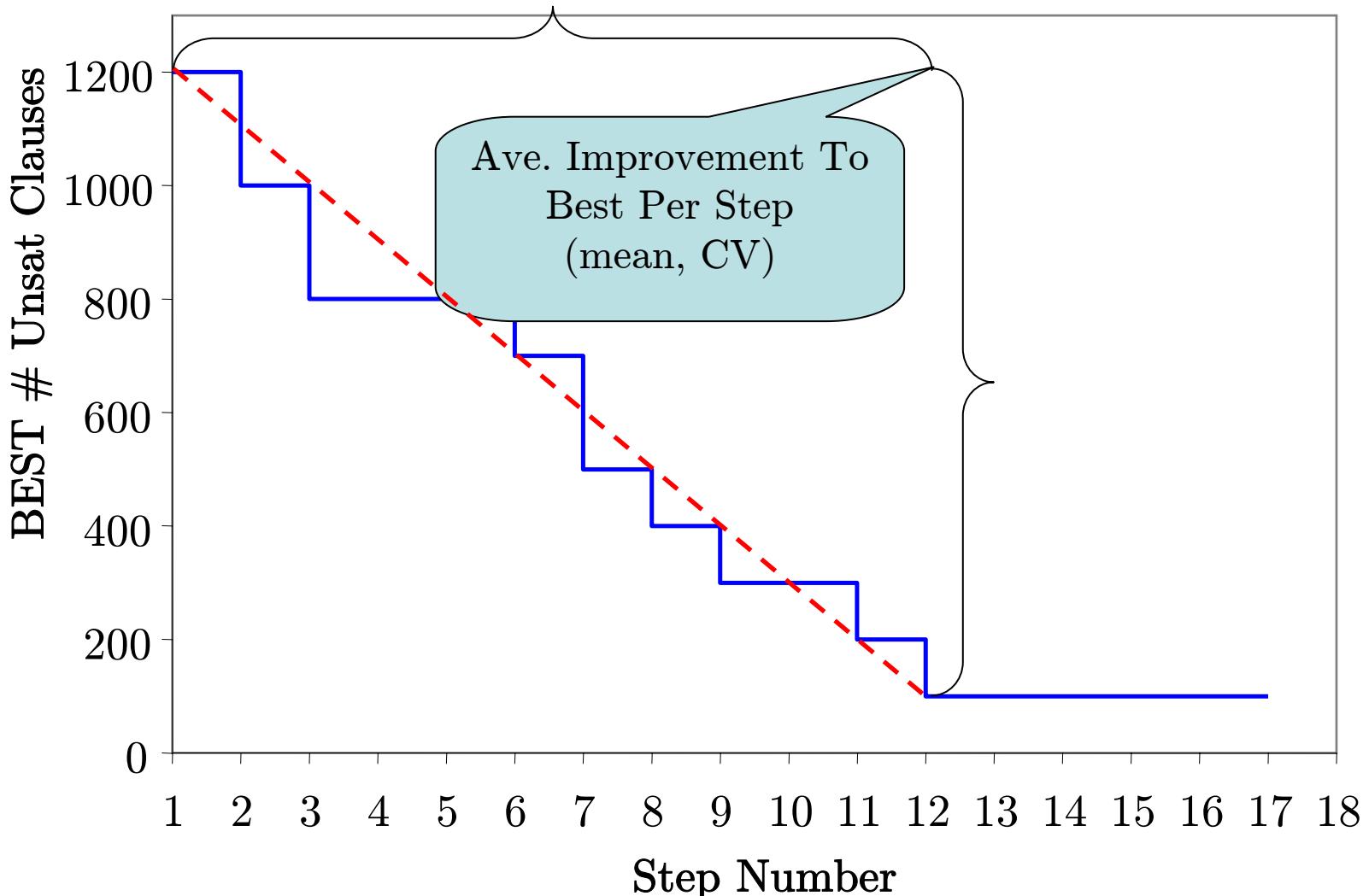
# Features: Local Search Probing



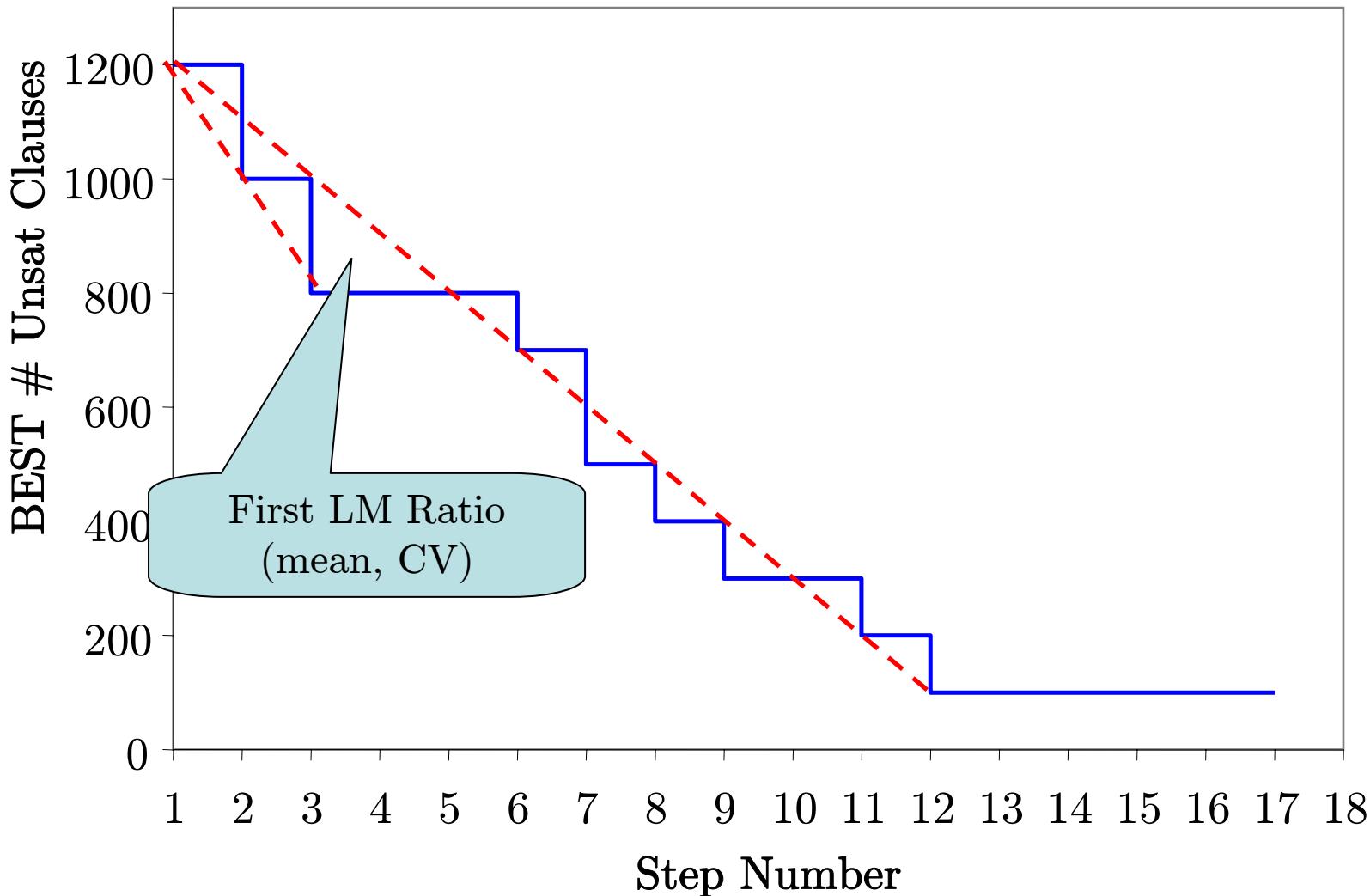
# Features: Local Search Probing



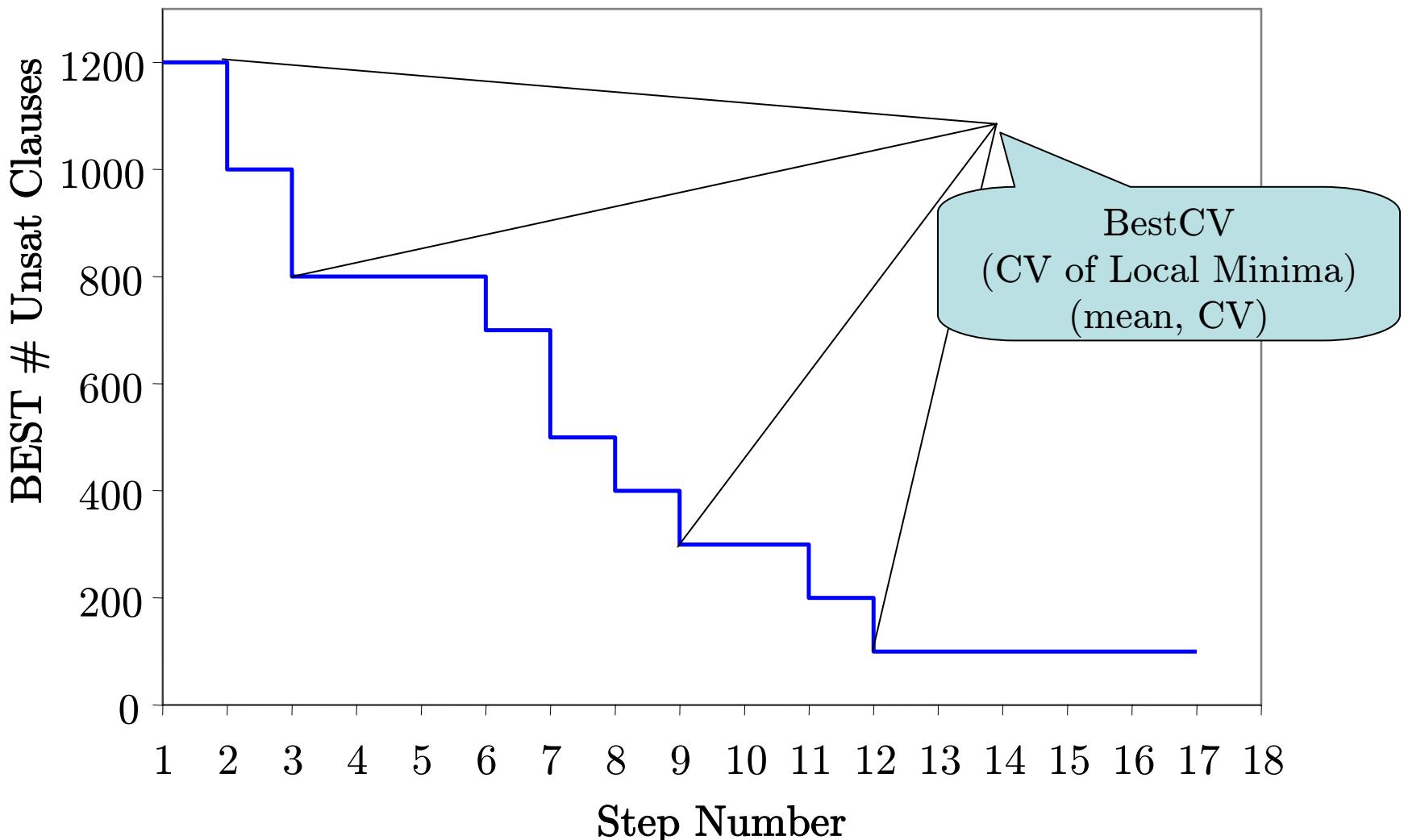
# Features: Local Search Probing



# Features: Local Search Probing

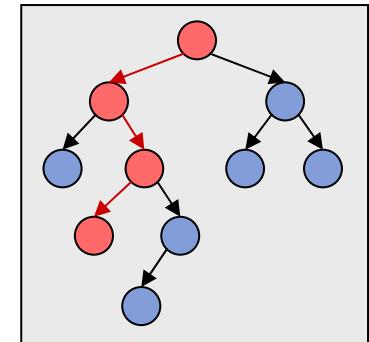


# Features: Local Search Probing



# Features: DPLL, LP

- **DPLL** search space size estimate
  - Random probing with unit propagation
  - Compute mean depth till contradiction
  - Estimate  $\log(\# \text{nodes})$



- Cumulative number of **unit propagations** at different depths (DPLL with Satz heuristic)

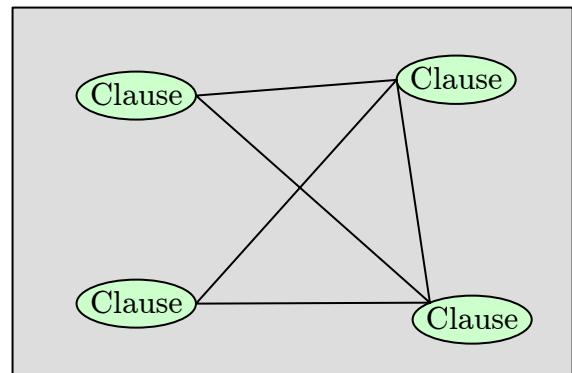
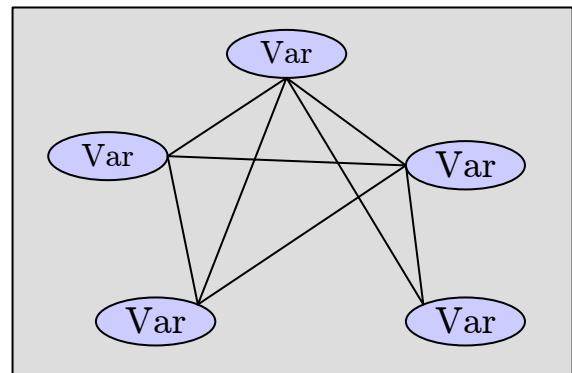
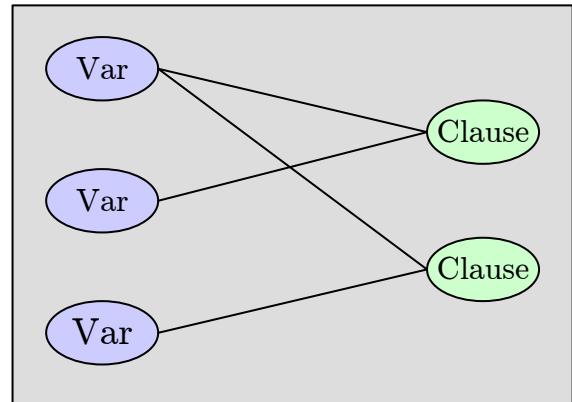
## • LP relaxation

- Objective value
- stats of integer slacks
- #vars set to an integer

$$\begin{aligned} \text{maximize: } & \sum_{k \in C} \left( \sum_{i \in L, i \in k} v_i + \sum_{j \in \bar{L}, i \in k} (1 - v_j) \right) \\ \text{subject to: } & \sum_{i \in k, i \in L} v_i + \sum_{j \in k, j \in \bar{L}} (1 - v_j) \geq 1 \quad \forall k \in C \\ & v_i \in \{0, 1\} \quad \forall i \end{aligned}$$

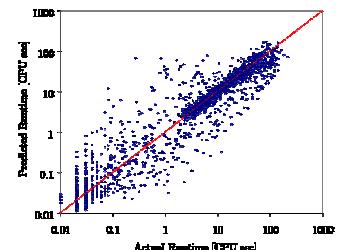
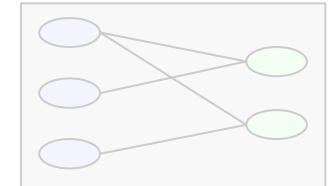
# Other Features

- **Problem Size:**
  - $v$  (#vars)
  - $c$  (#clauses)
  - Powers of  $c/v, v/c, |c/v - 4.26|$
- **Graphs:**
  - **Variable-Clause** (VCG, bipartite)
  - **Variable** (VG, edge whenever two variables occur in the same clause)
  - **Clause** (CG, edge iff two clauses share a variable with opposite sign)
- **Balance**
  - #pos vs. #neg literals
  - unary, binary, ternary clauses
- Proximity to **Horn formula**



# Outline

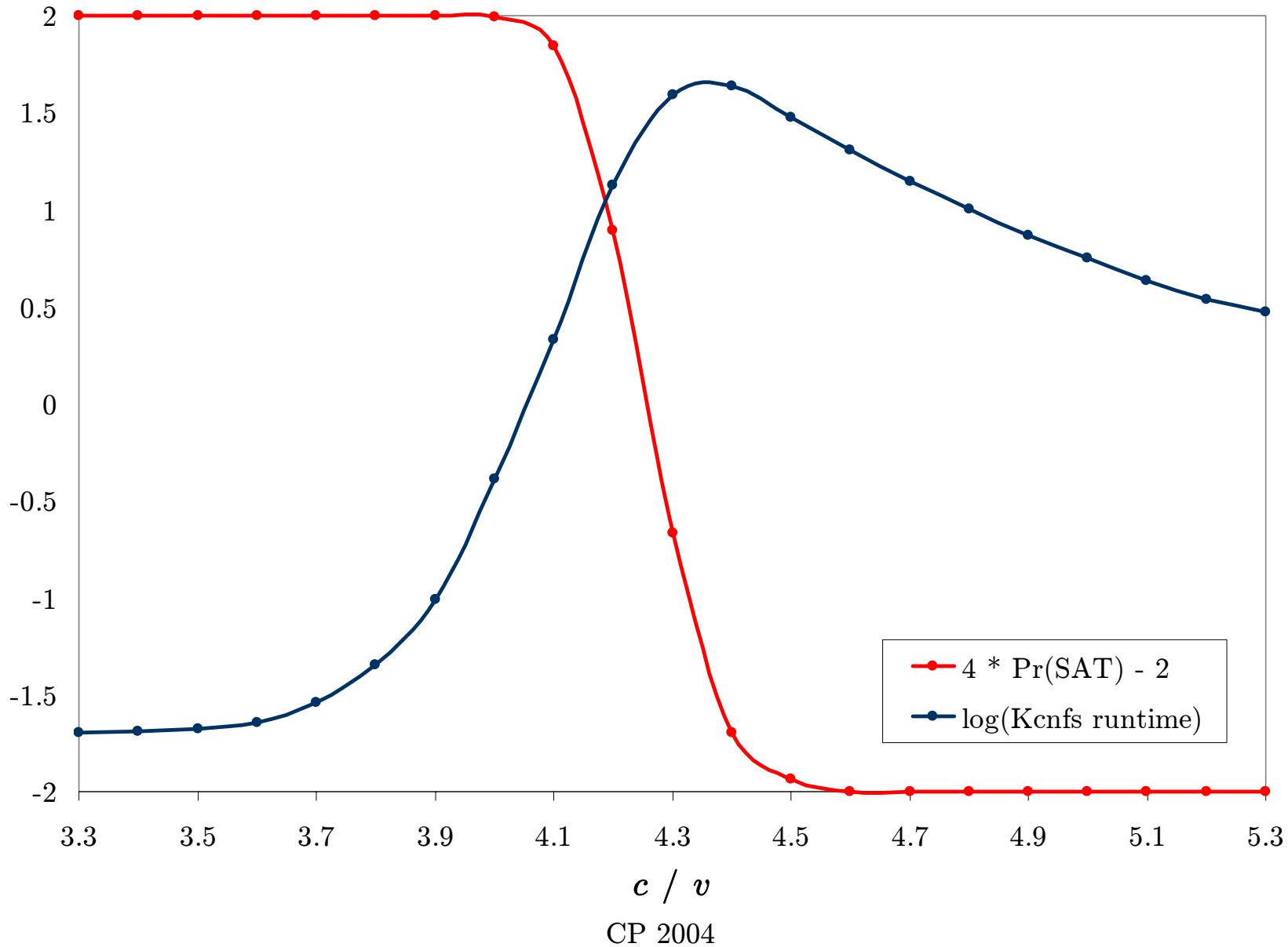
- Features
- Experimental Results
  - Variable Size Data
  - Fixed Size Data



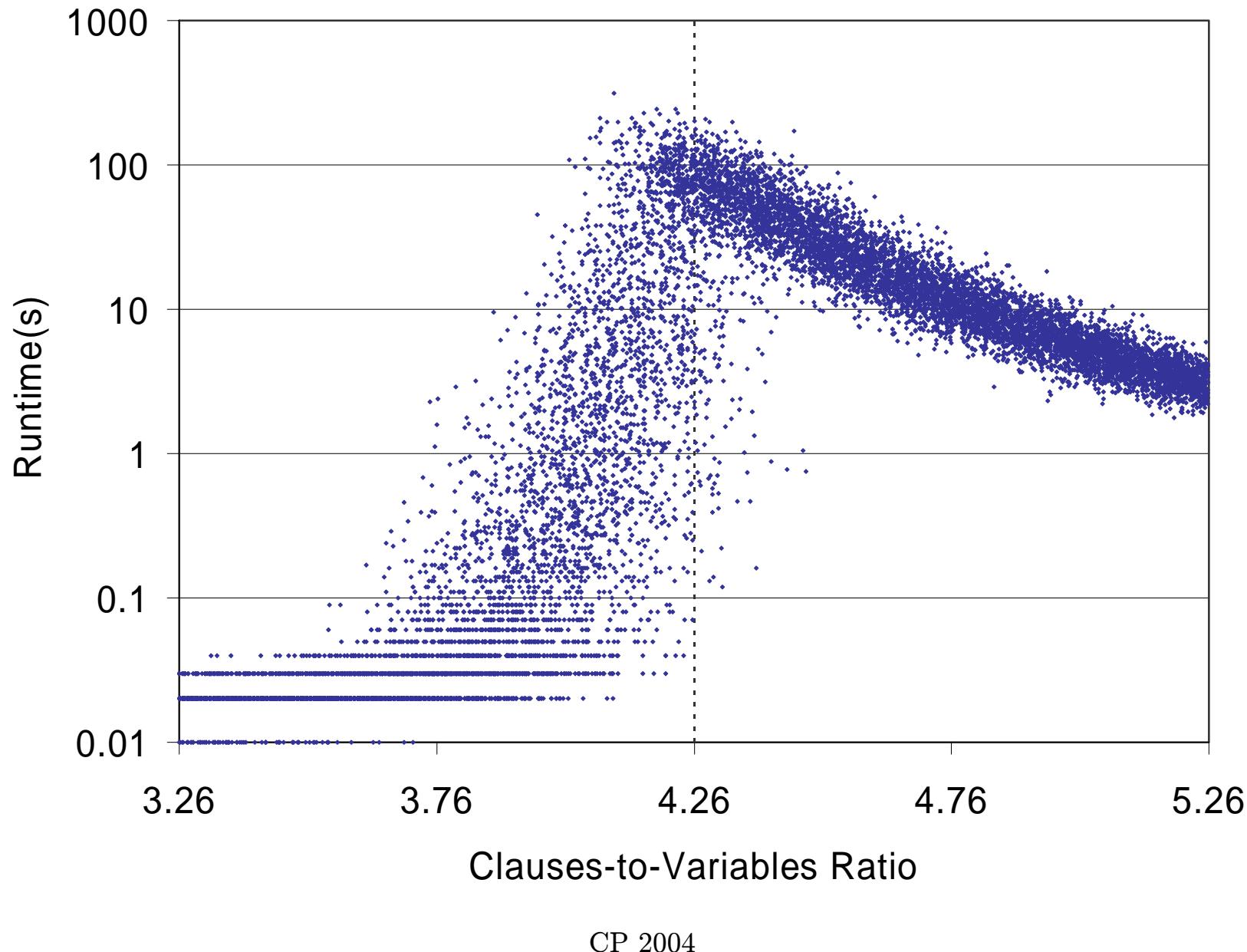
# Experimental Setup

- Uniform random 3-SAT, 400 vars
- **Datasets** (20000 instances each)
  - **Variable-ratio** dataset (1 CPU-month)
    - $c/v$  uniform in [3.26, 5.26] ( $\therefore c \in [1304, 2104]$ )
  - **Fixed-ratio** dataset (4 CPU-months)
    - $c/v=4.26$  ( $\therefore v=400, c=1704$ )
- **Solvers**
  - Kcnfs [Dubois and Dequen]
  - OKsolver [Kullmann]
  - Satz [Chu Min Li]
- **Quadratic regression** with logistic response function
- Training : test : validation split – 70 : 15 : 15

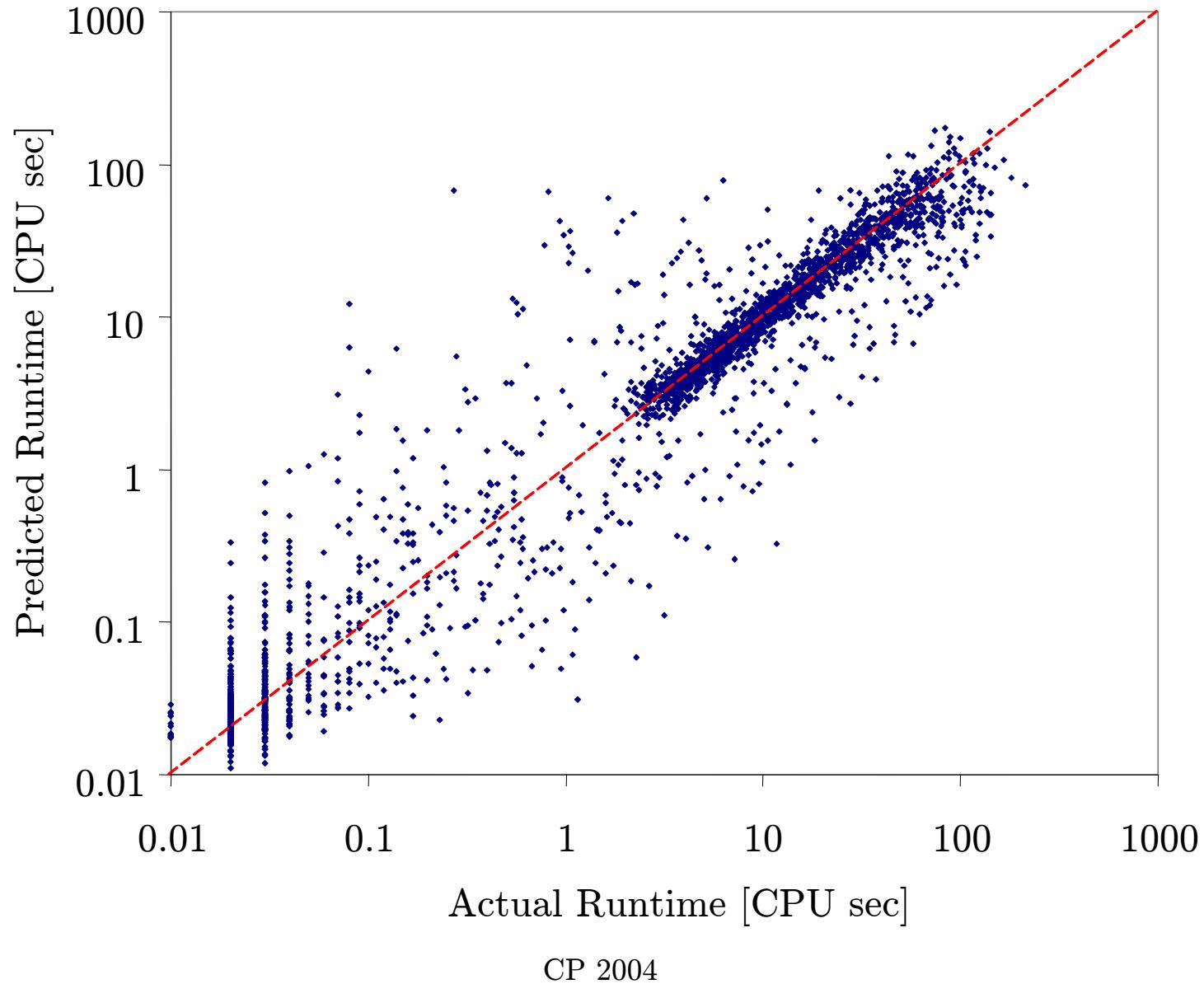
# Kcnfs Data



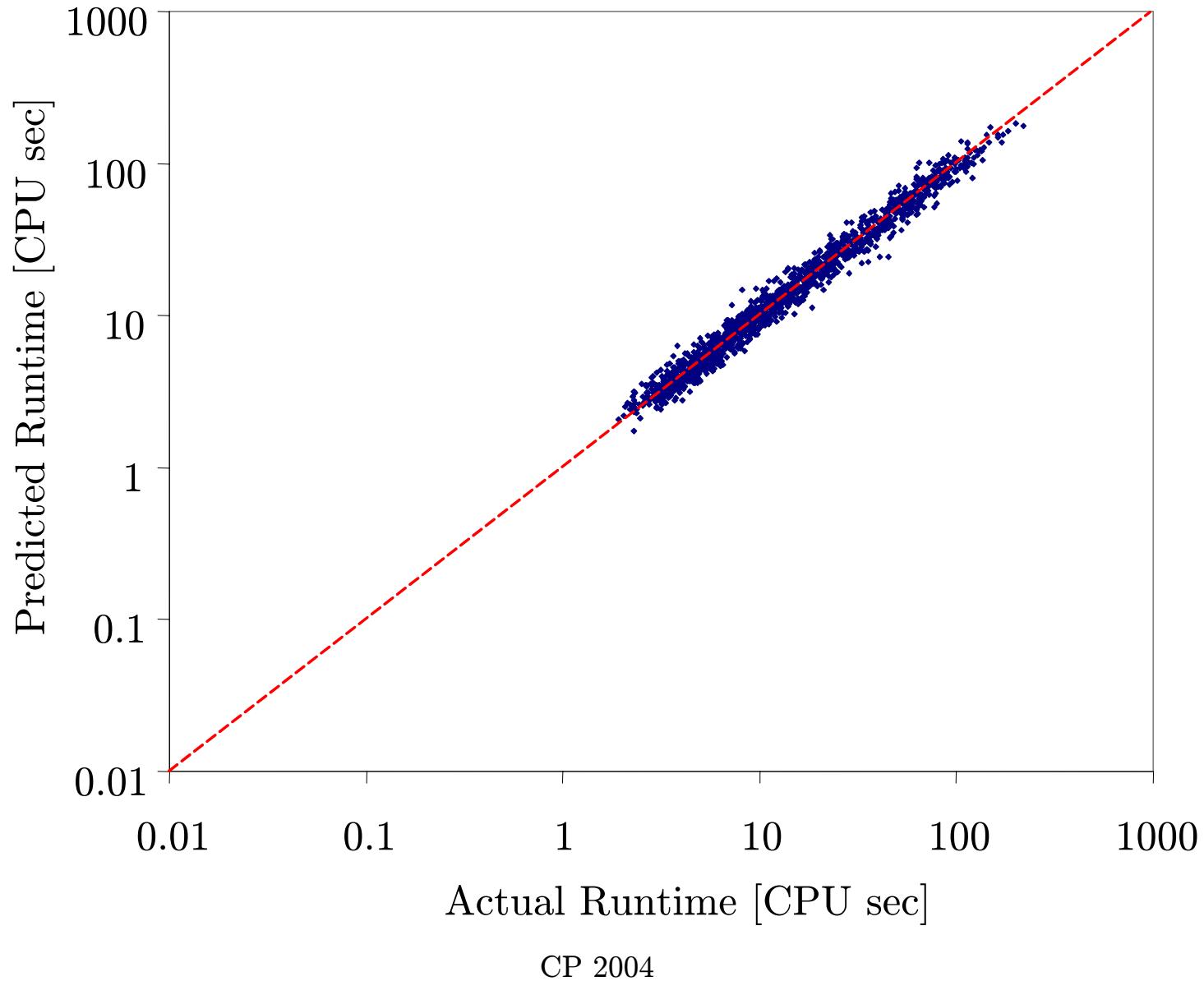
# Kcnfs Data



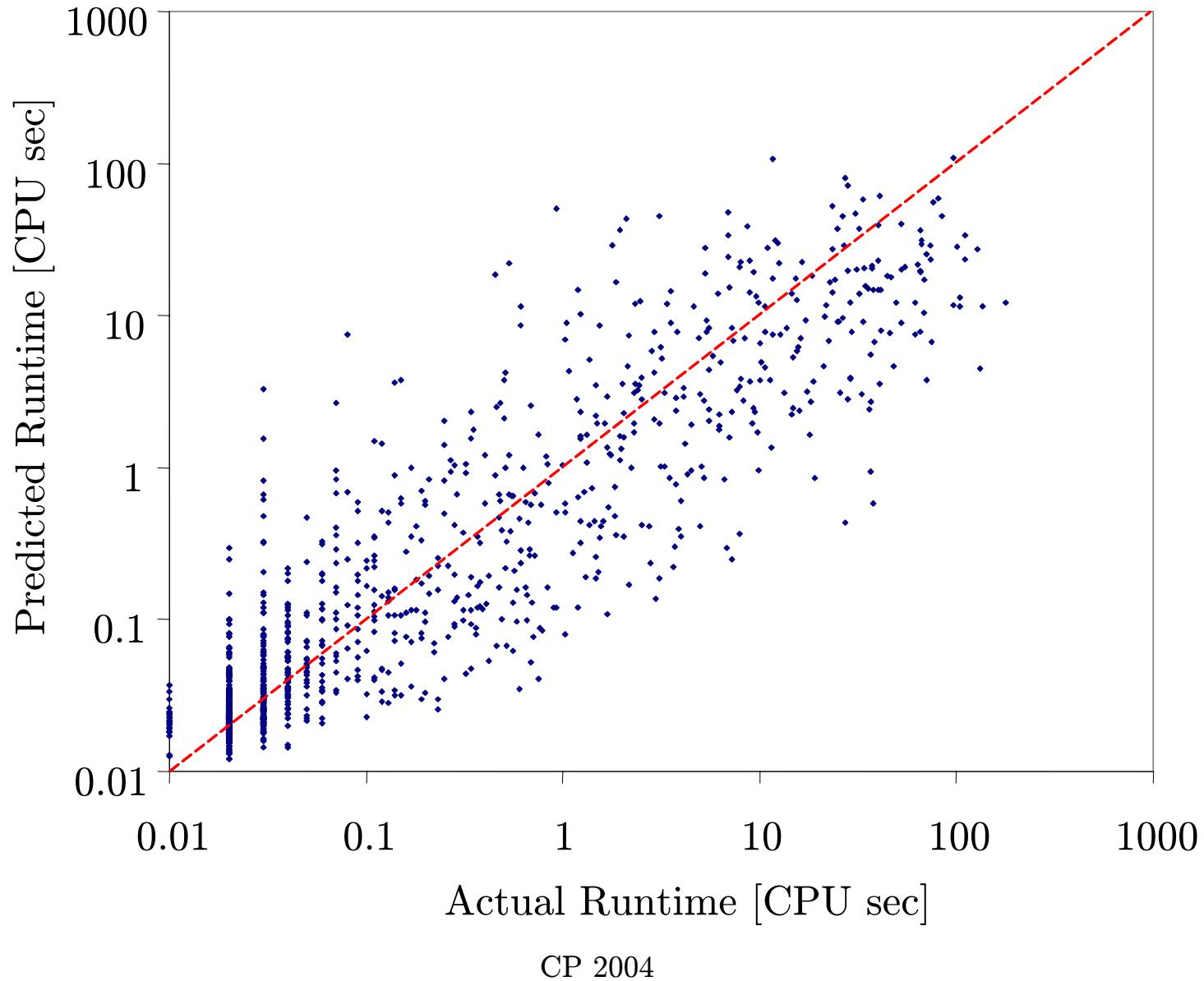
# Variable Ratio Prediction (Kcnfs)



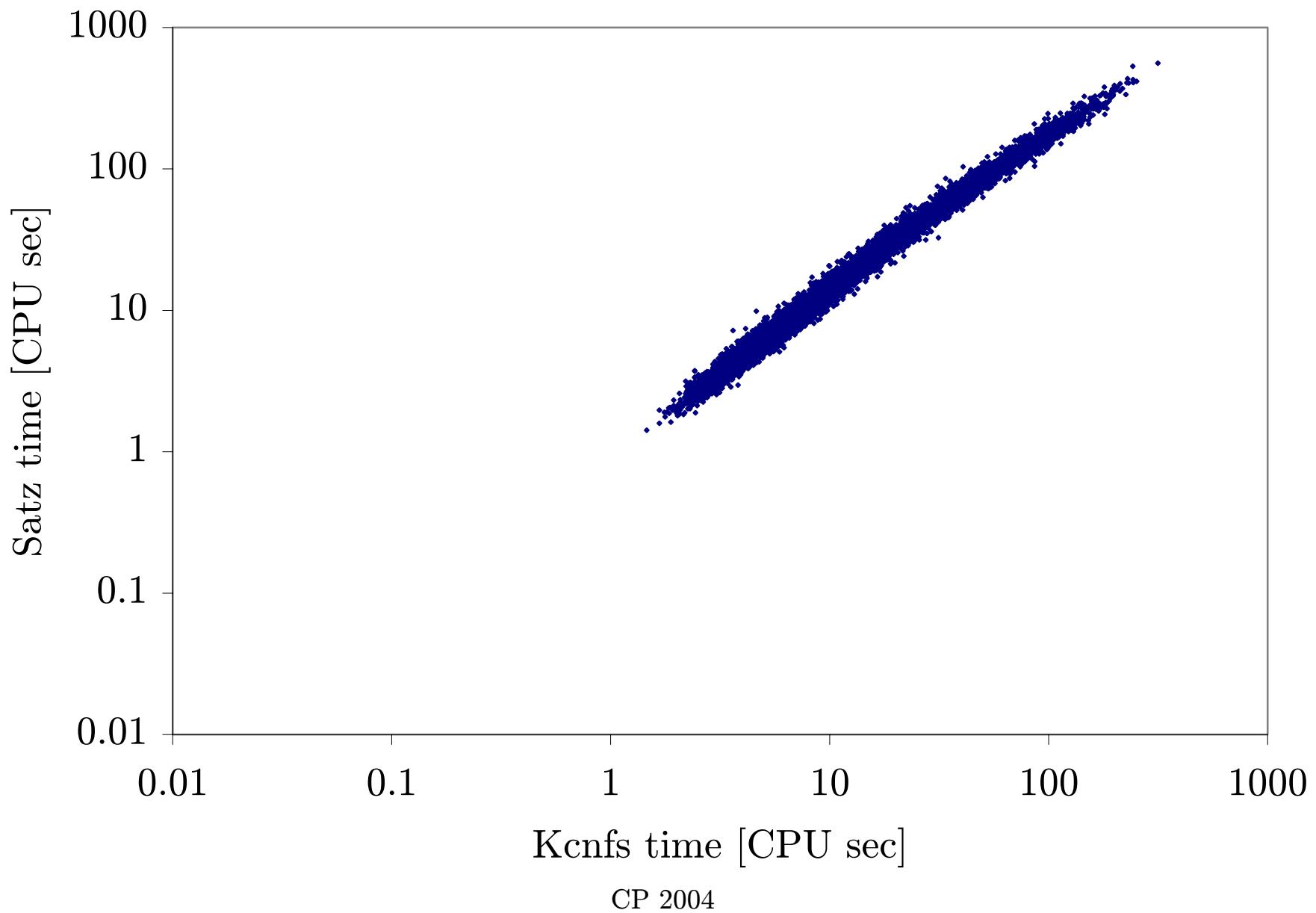
# Variable Ratio - UNSAT



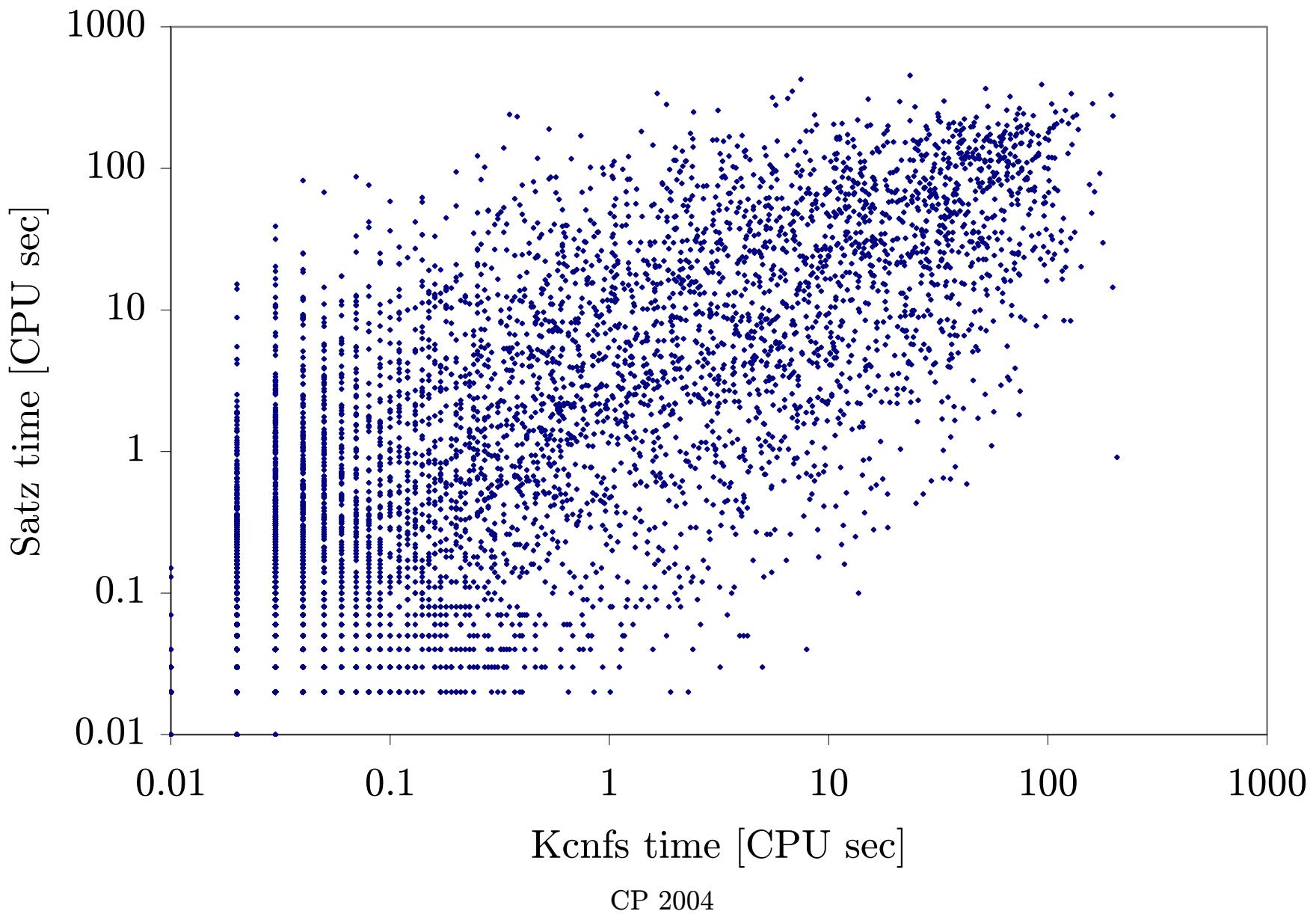
# Variable Ratio - SAT



# Kcnfs vs. Satz (UNSAT)



# Kcnfs vs. Satz (SAT)

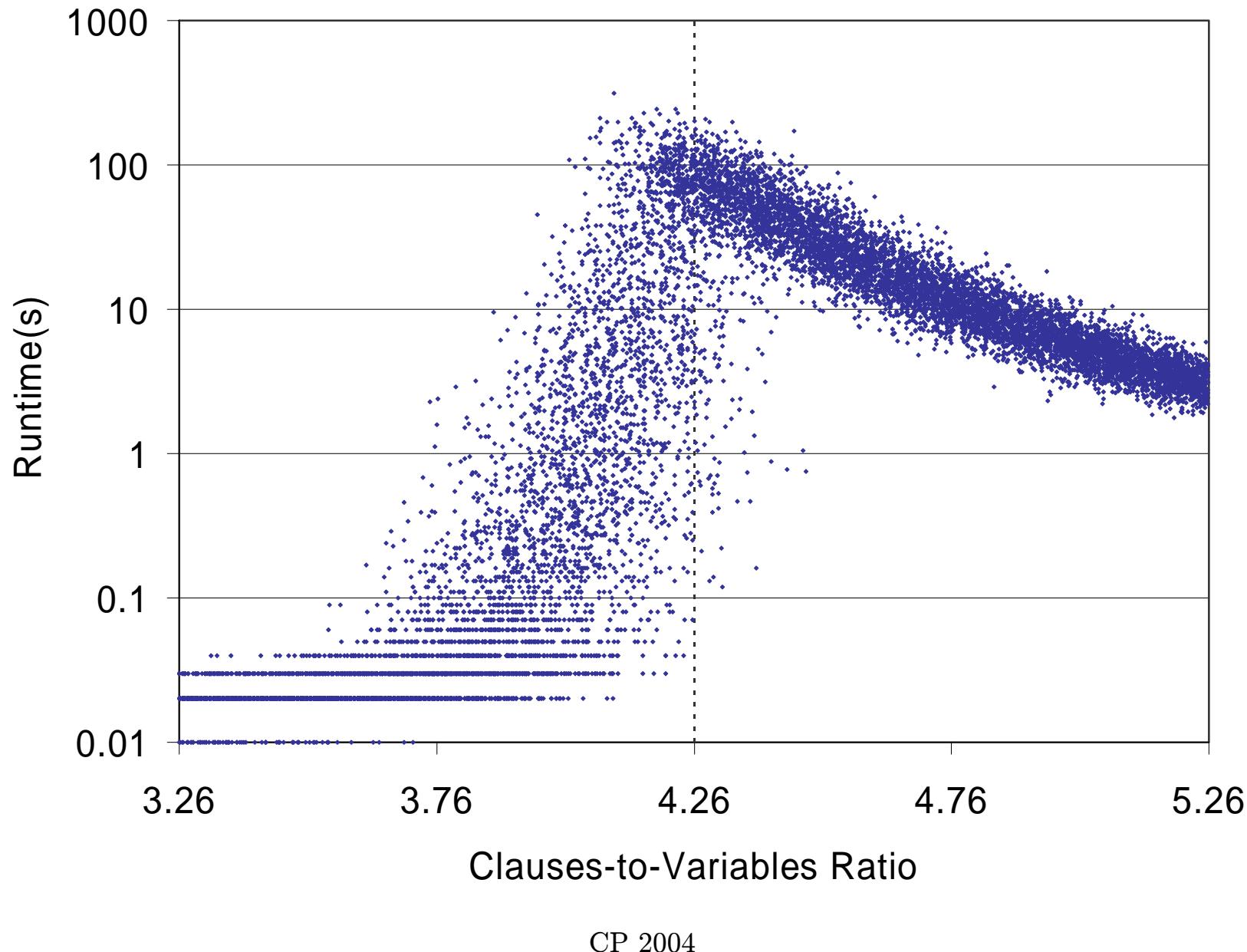


# Feature Importance – Variable Ratio

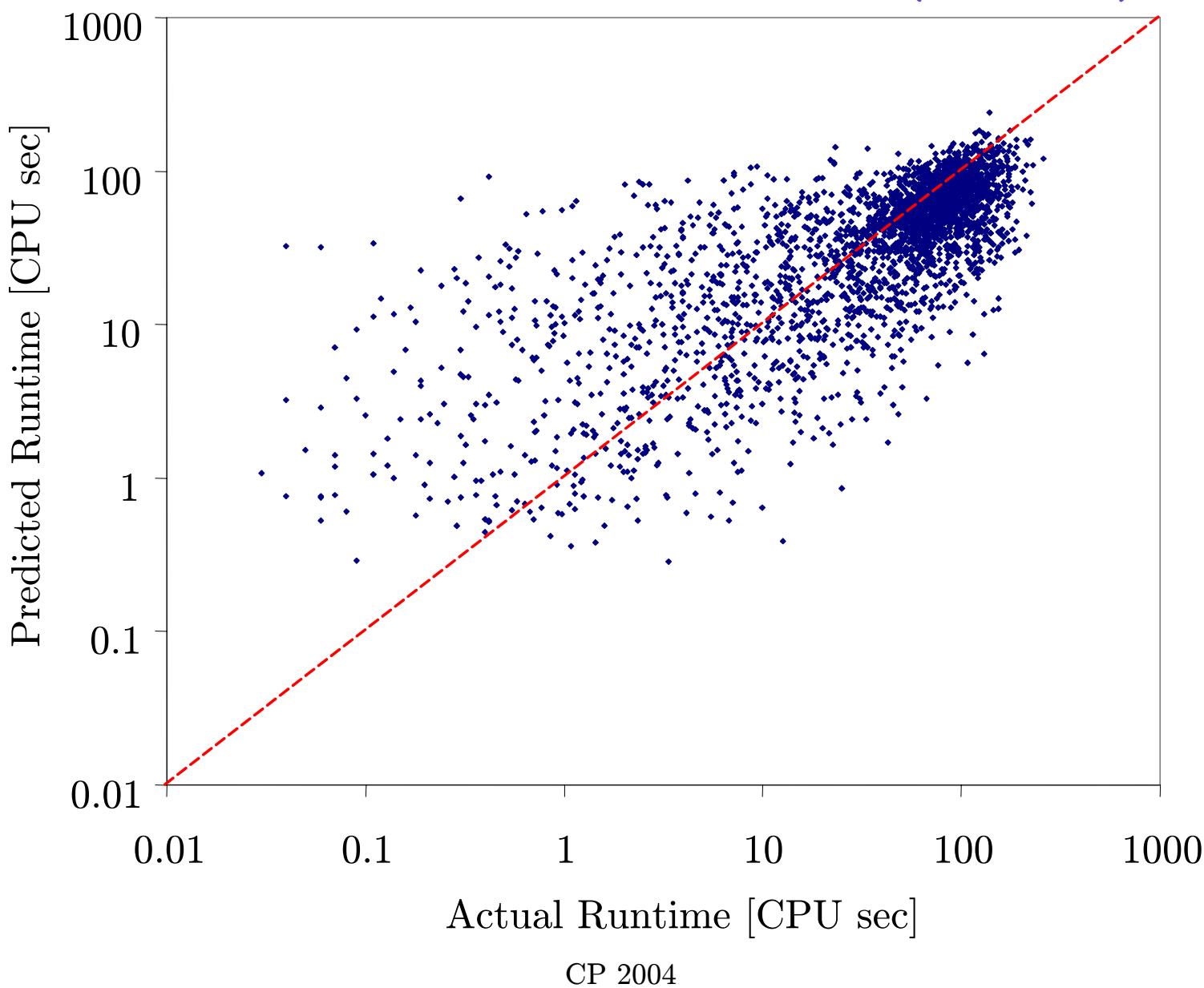
- **Subset selection** can be used to identify features **sufficient** for approximating full model performance
- Other (correlated) sets could potentially achieve similar performance

Variable	Cost of Omission
$ c/v-4.26 $	100
$ c/v-4.26 ^2$	69
$(v/c)^2 \times \text{SapsBestCVMean}$	53
$ c/v-4.26  \times \text{SapsBestCVMean}$	33

# Fixed Ratio Data



# Fixed Ratio Prediction (Kcnfs)



# Feature Importance – Fixed Ratio

Variable	Cost of Omission
$SapsBestSolMean^2$	100
$SapsBestSolMean \times \text{MeanDPLLDepth}$	74
$GsatBestSolCV \times \text{MeanDPLLDepth}$	21
$VCGClauseMean \times GsatFirstLMRatioMean$	9

*(Note: The first four rows correspond to the variables shown in the table above.)*

**Diagram 1:** A directed acyclic graph (DAG) illustrating variable dependencies. Nodes are colored red or blue. Red nodes represent variables, and blue nodes represent clauses. Arrows indicate dependencies from variables to clauses. The graph shows a complex dependency structure where multiple variables influence many clauses.

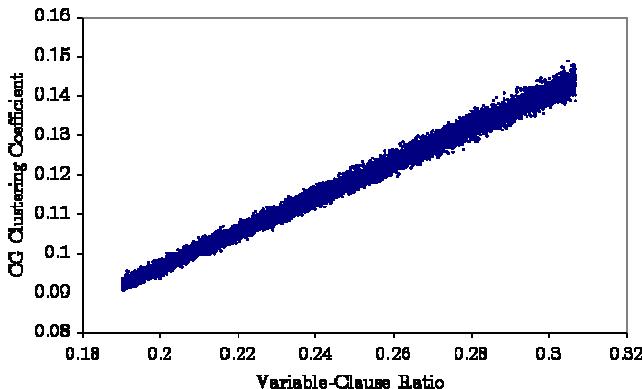
**Diagram 2:** A bipartite graph illustrating the relationship between variables and clauses. Nodes are colored purple (Variables) and green (Clauses). Edges connect variables to clauses, showing that each variable can be associated with multiple clauses and vice versa.

# SAT vs. UNSAT

- Training models separately for SAT and UNSAT instances:
  - good models require fewer features
  - model accuracy improves
  - $c/v$  no longer an important feature for VR data
  - Completely different features are useful for SAT than for UNSAT
- Feature importance on SAT instances:
  - Local Search features sufficient
    - 7 features for good VR model
    - 1 feature for good FR model ( $SAPSBestSolCV \times SAPSAveImpMean$ )
  - If LS features omitted, LP + DPLL search space probing
- Feature importance on UNSAT instances:
  - DPLL search space probing
  - Clause graph features

# Beyond Ratio: Weighted CG Clustering Coefficient

- Byproduct of our analysis: a very **strong correlation** between weighted CG clustering coefficient and  $v/c$



- Clustering coefficient is a more fundamental concept than  $v/c$ , since it describes the **structure of the constraints** explicitly, not implicitly.
  - correlation between (unweighted) CC and hardness has been shown for **other constraint problems** (e.g., graph coloring, combinatorial auctions)
- We have a **proof sketch** of this correlation

# Conclusions

- Can construct **good models** for DPLL solvers
- These models can be **analyzed** to gain understanding about what makes instances hard or easy for solvers
- **Algorithm portfolios** can be constructed (**Satzilla**)
- More specifically:
  - Strong relationship between **LS** and **DPLL** search spaces
  - Our approach **automatically identified** importance of  $c/v$
  - **SAT/UNSAT instances** have very different performance characteristics; it helps to model them separately
  - **Clustering Coefficient** explains why  $c/v$  is important in terms of local properties of constraint graph