# Προβλήματα Ικανοποίησης Περιορισμών: από τη Φυσική στους Αλγορίθμους
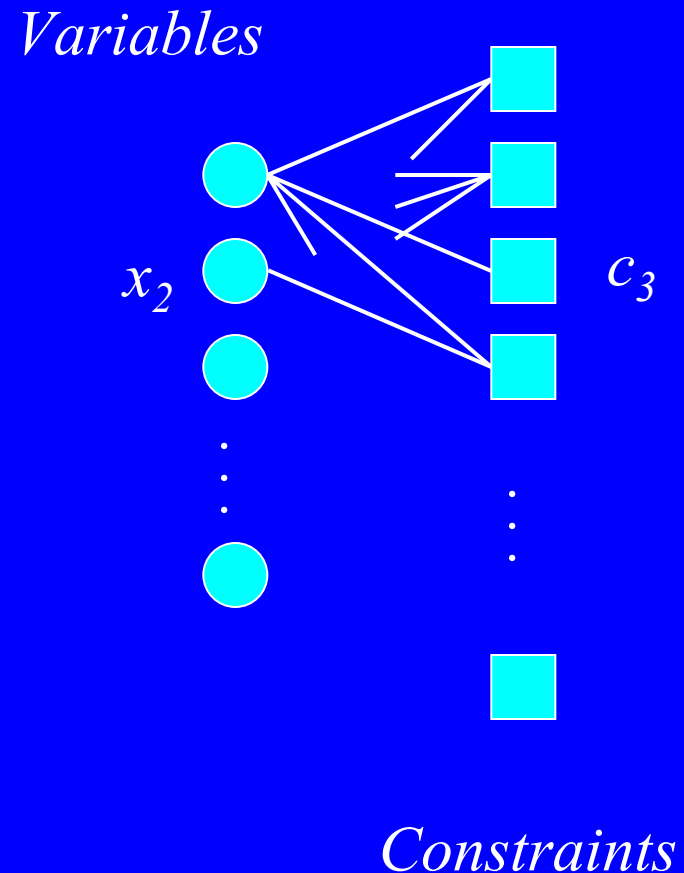
Δημήτρης Αχλιόπτας

University of California
Santa Cruz

# The Setting: Random CSPs

- n variables with small, discrete domains
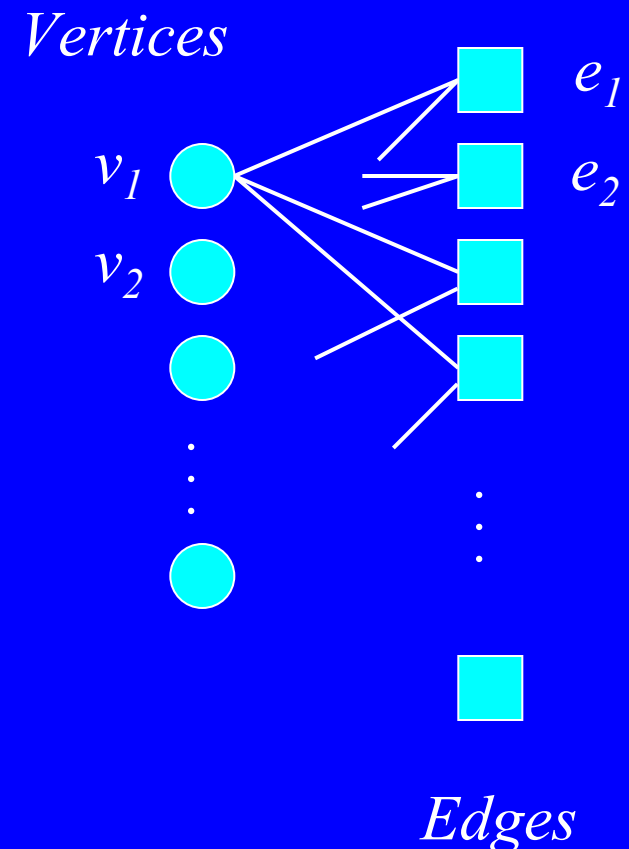
- m conflicting constraints

---

- Random bipartite graph:

- Sparse graph, i.e. m=Θ(n)

*Variables*

$x_2$

$c_3$

*Constraints*

# Random Graph k-coloring

- Each vertex is a variable with domain $\{1,2,\ldots,k\}$
- Each edge is a "not-equal" constraint on two variables

---

- G(n,m) random graph: the two variables are chosen randomly
- Random r-regular: each variable is chosen r times

*Vertices*

$v_1$

$v_2$

$e_1$

$e_2$

*Edges*

# Random k-SAT

- Take $n$ Boolean variables $X = \{x_1, x_2, \ldots, x_n\}$

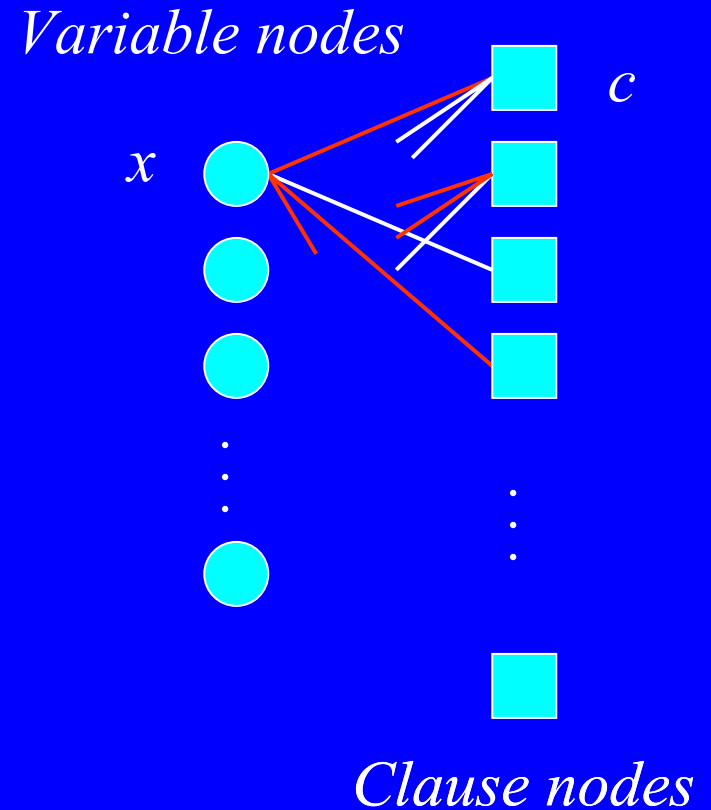- Among all $2^k \binom{n}{k}$ possible k-clauses select $m$

  uniformly and independently. Typically $m = rn$ .

- Example ( $k = 3$ ) :
$(\overline{x}_{12} \vee x_5 \vee \overline{x}_9) \wedge (x_{34} \vee \overline{x}_{21} \vee x_5) \wedge \cdots \cdots \wedge (x_{21} \vee x_9 \vee \overline{x}_{13})$

# Random k-SAT

- Variables are binary.
- Every constraint (k-clause) binds k variables.
- Forbids exactly one of the $2^k$ possible joint values.

- Random k-SAT = each clause picks k random literals.

*Variable nodes*

*x*

*c*

*Clause nodes*

# Two Values

**Theorem.** For every $d > 0$, w.h.p. the chromatic number of $G(n, p = d/n)$

$$\text{is either } k \text{ or } k + 1$$

where $k$ is the smallest integer s.t. $d < 2k \log k$.

[A., Naor '04]

# Examples

- If $d = 7$, w.h.p. the chromatic number is $4$ or $5$ .

- If $d = 10^{60}$, w.h.p. the chromatic number is

$$3771455490672260758090142394938336005516126417647650681575$$

or

$$3771455490672260758090142394938336005516126417647650681576$$

# A simple k-coloring algorithm

- Repeat
  - Pick a random uncolored vertex
  - Assign it the lowest allowed number (color)

Works when $d \leq k \log k$ [Bollobás, Thomasson 84]
[McDiarmid 84]

- There are no $k$-colorings for $d \geq 2k \log k$

# The satisfiability threshold conjecture

Conjecture: for every $k \geq 3$, there is $r_k$ such that

$$\lim_{n \to \infty} \Pr[\mathcal{F}_k(n, rn) \text{ is satisfiable}] = \begin{cases} 1 & \text{if } r = r_k - \epsilon \\ 0 & \text{if } r = r_k + \epsilon \end{cases}$$

Since the 80s: for every $k \geq 3$,

$$c \; \frac{2^k}{k} < r_k < 2^k \ln 2$$

[Chvátal & Reed 92]

[Frieze & Suen 96]

# Easy Upper Bound

The probability there is a satisfying assignments is at most:

$$2^n \left(1 - \frac{1}{2^k}\right)^m = \left[2\left(1 - \frac{1}{2^k}\right)^r\right]^n$$

$$\to 0 \quad \text{for } r \geq 2^k \ln 2$$

# Lower Bound

Repeat:
- Pick a random variable and set it randomly
- Satisfy 1-clauses if they exist (repeatedly)
- Fail if any 0-clause occurs

● Finds a satisfying truth assignment w.h.p. for all

$$r < \frac{2^k}{k}$$    [Chao & Franco '86]

# Bounds for the k-SAT threshold

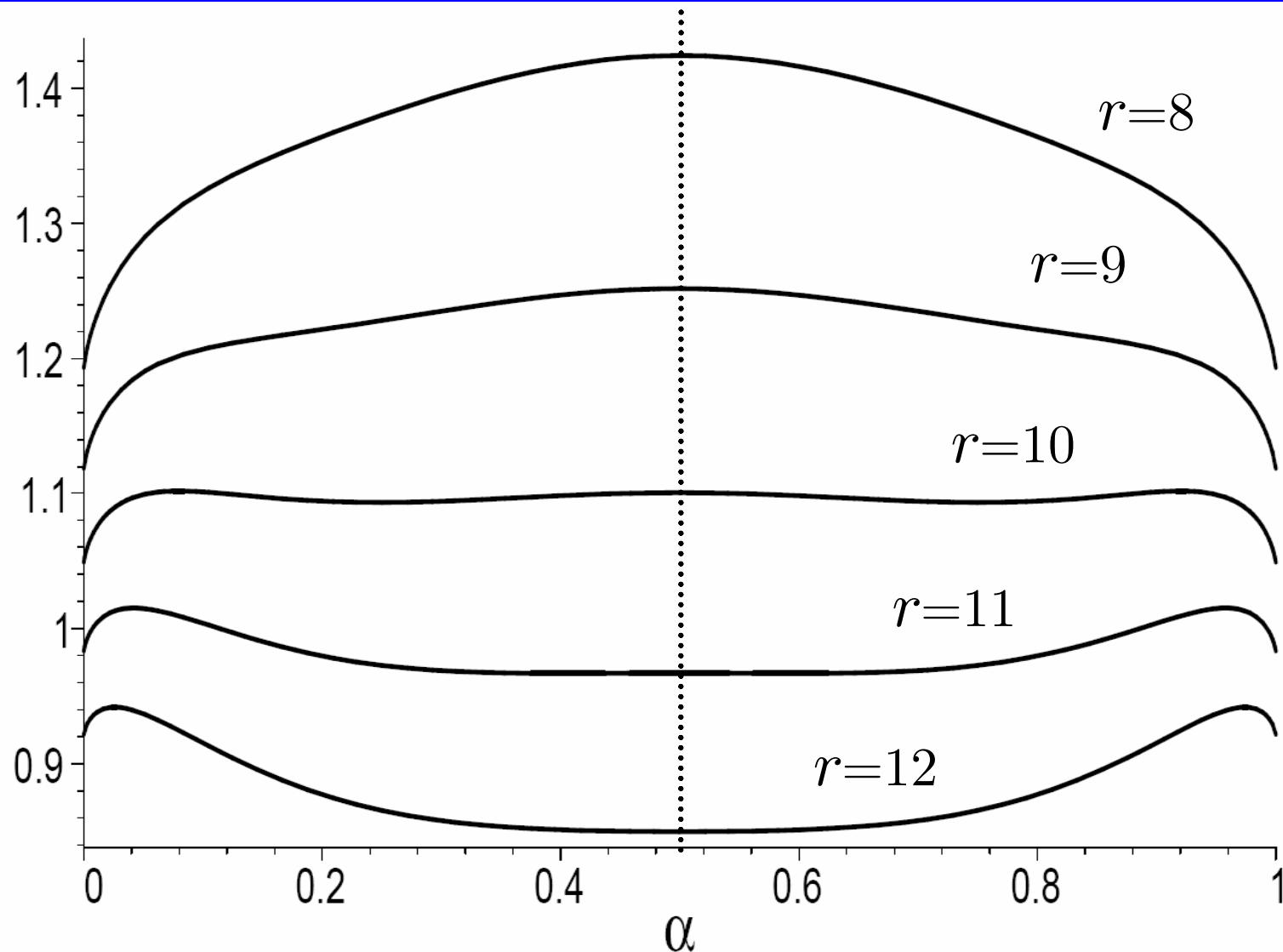For all $k \geq 3$:

$$2^k \ln 2 - k < r_k < 2^k \ln 2$$

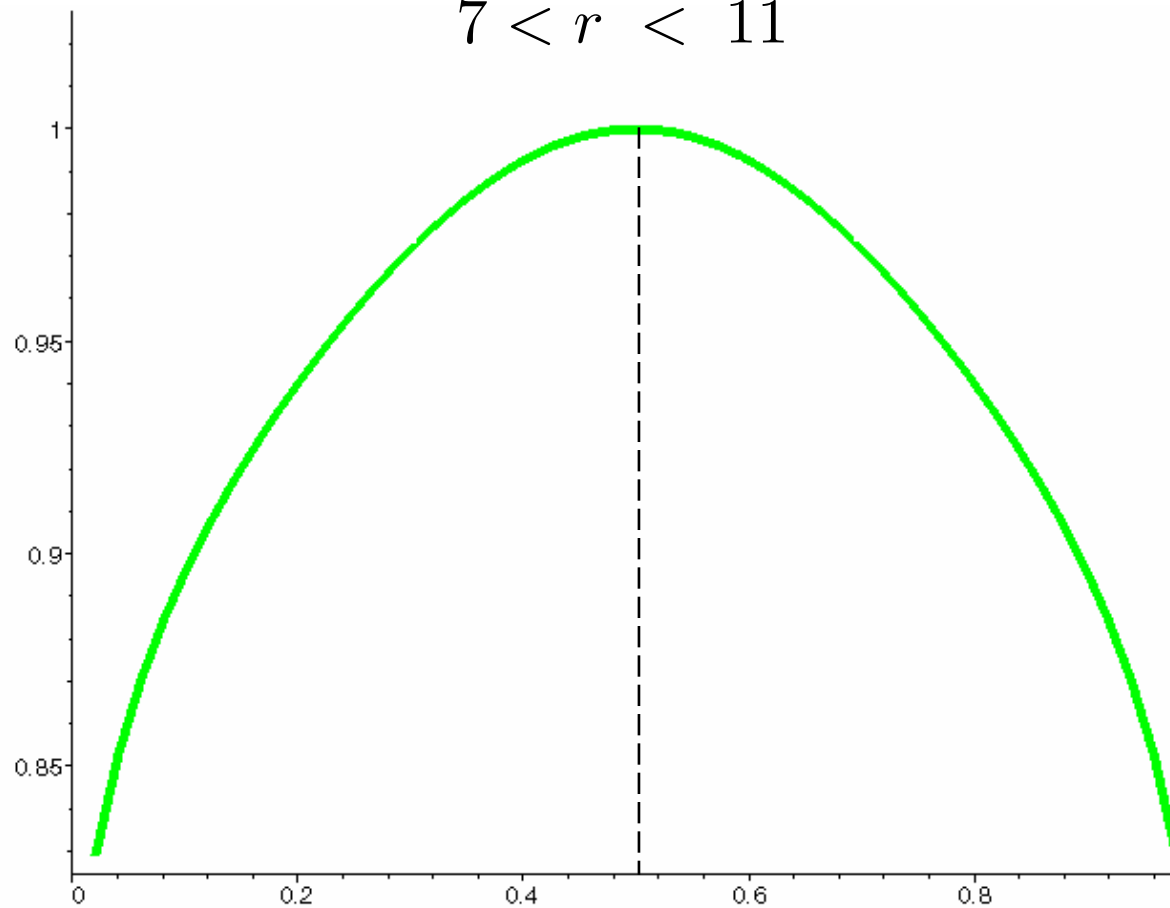| $k$ | 3 | 4 | 5 | 7 | 10 | 20 | 21 |
|---|---|---|---|---|---|---|---|
| Upper bound | 4.51 | 10.23 | 21.33 | 87.88 | 708.94 | 726,817 | 1,453,635 |
| Lower bound | 3.52 | 7.91 | 18.79 | 84.82 | 704.94 | 726,809 | 1,453,626 |
| Best algorithm | 3.52 | 5.54 | 9.63 | 33.23 | 172.65 | 95,263 | 181,453 |

# Bicoloring 5-uniform hypergraphs

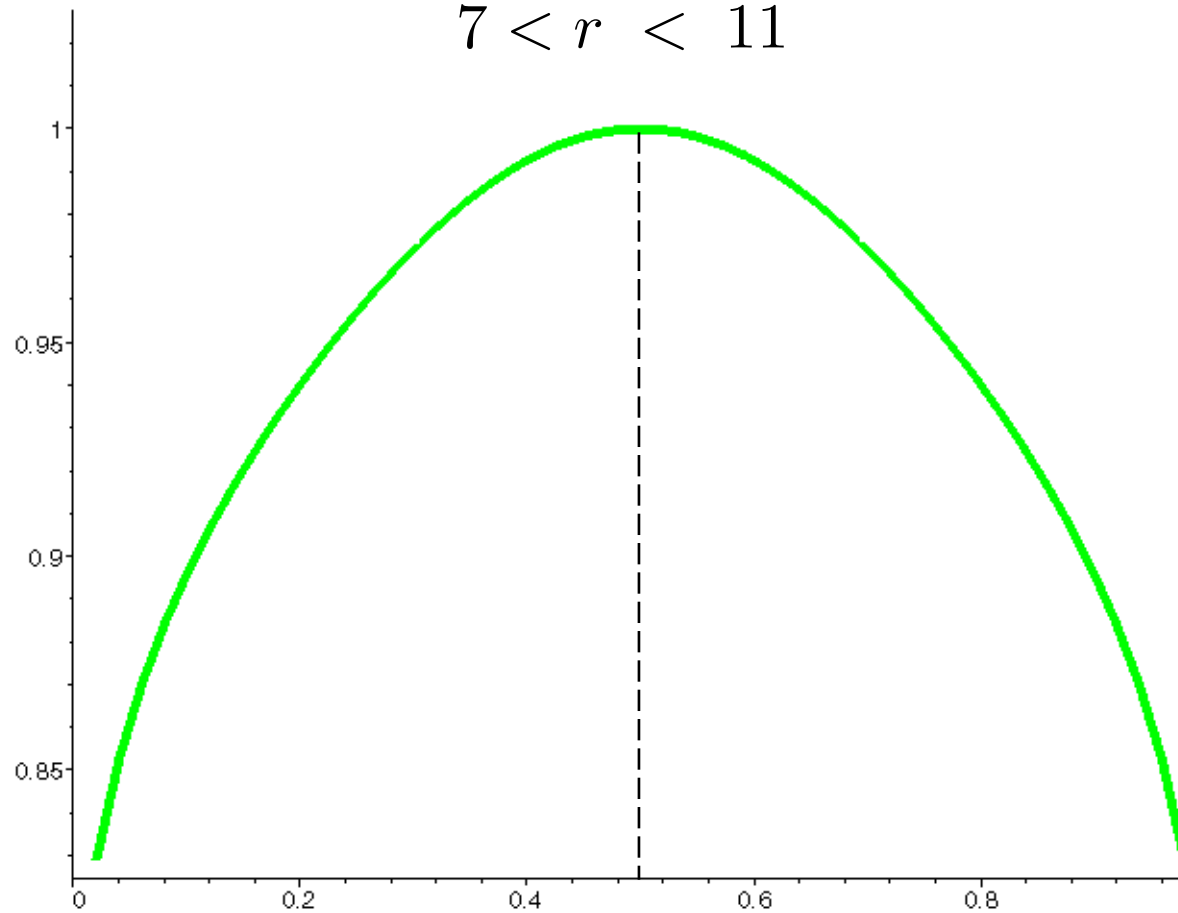# 5-uniform hypergraphs

$$7 < r < 11$$

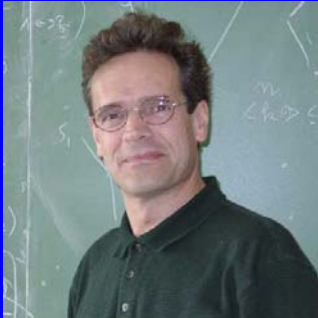# 5-uniform hypergraphs

$$7 < r < 11$$

# Natural question

Are there efficient algorithms
that work closer
to each problem's threshold?

# Our Best Algorithms are Naive

- Repeat
  - Pick a random uncolored vertex
  - Assign it the lowest available color

- Repeat
  - Pick a random variable and set it randomly
  - Satisfy 1-clauses if they exist (repeatedly)

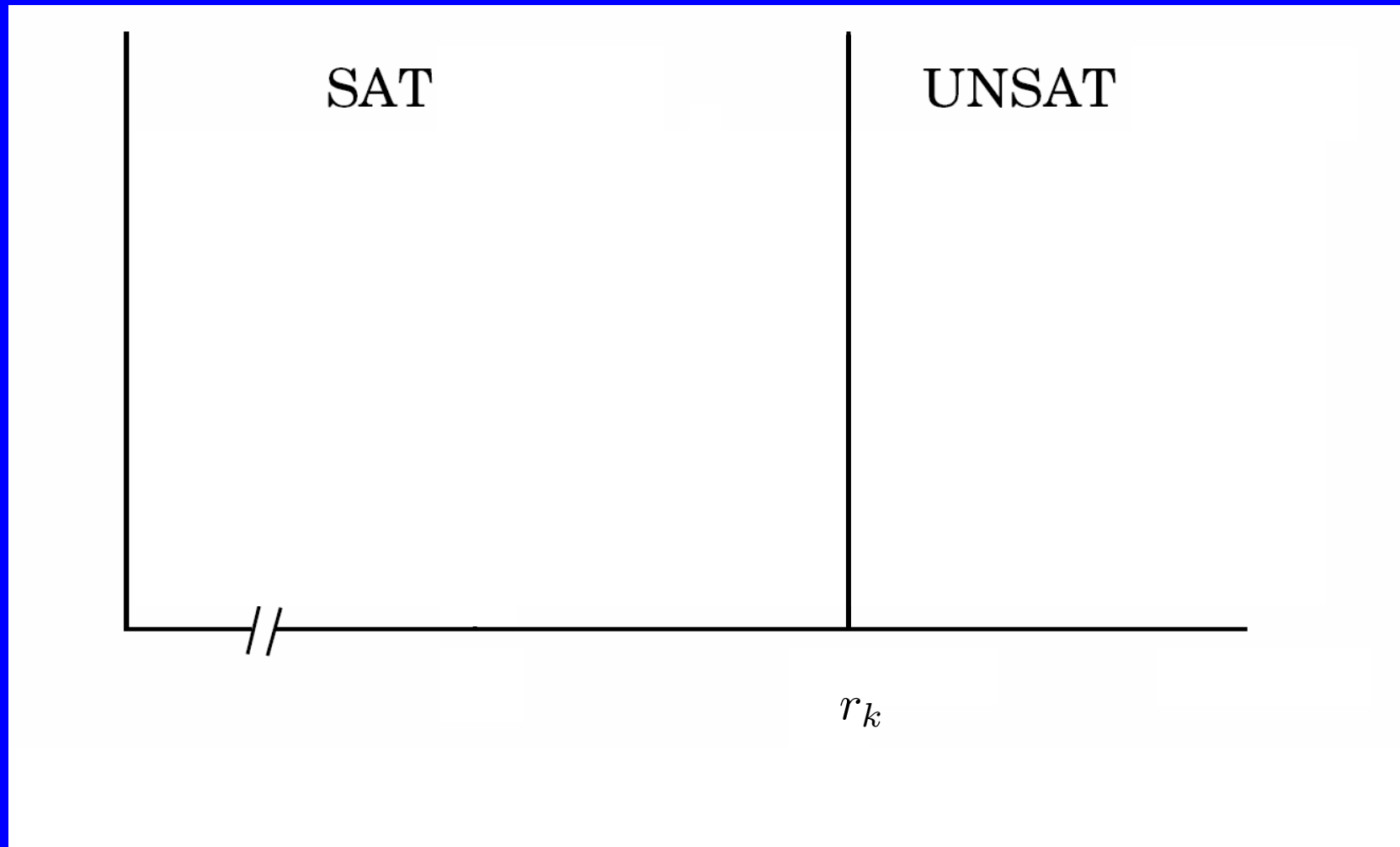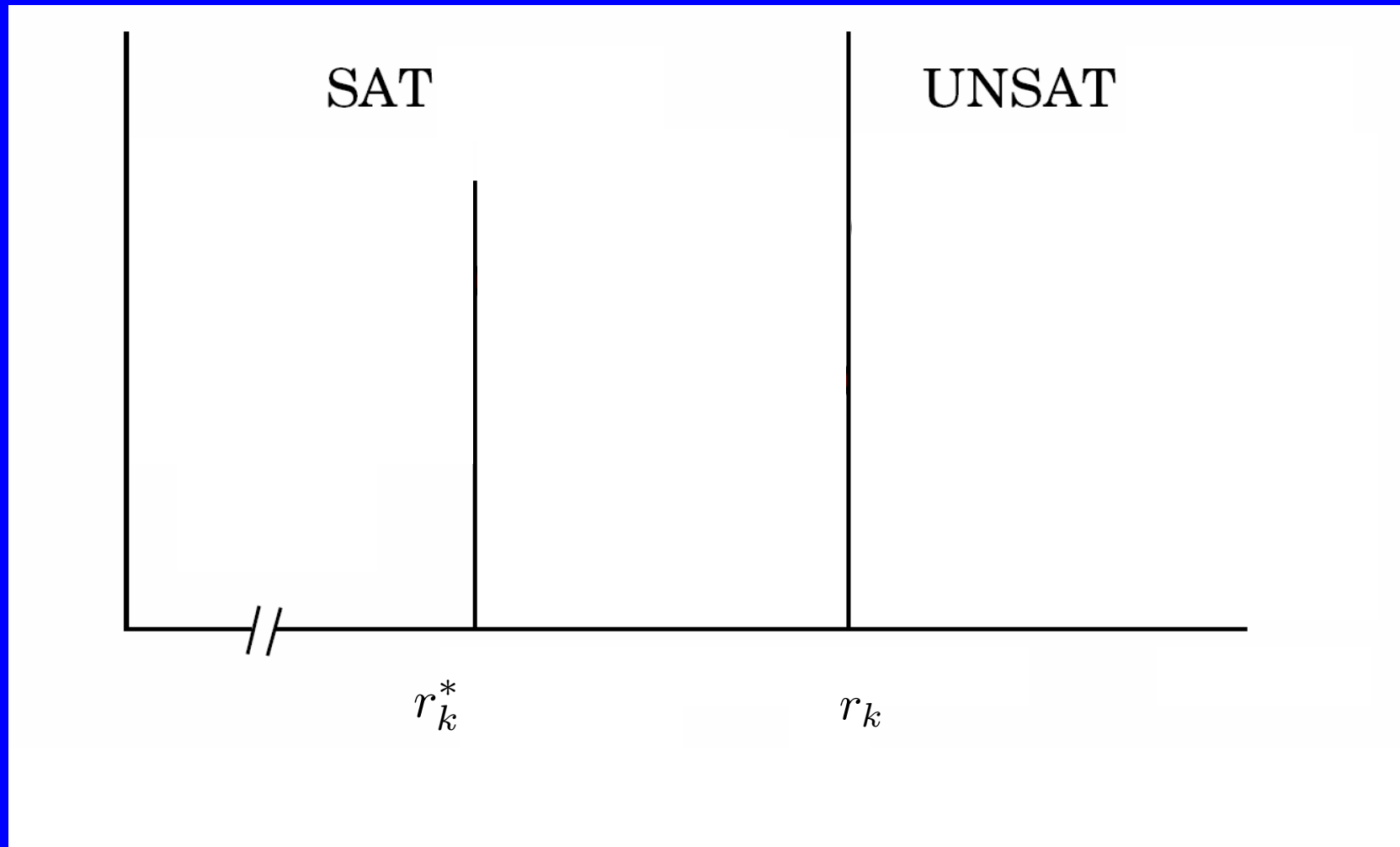# In a parallel universe



Marc Mézard     Giorgio Parisi     Riccardo Zecchina
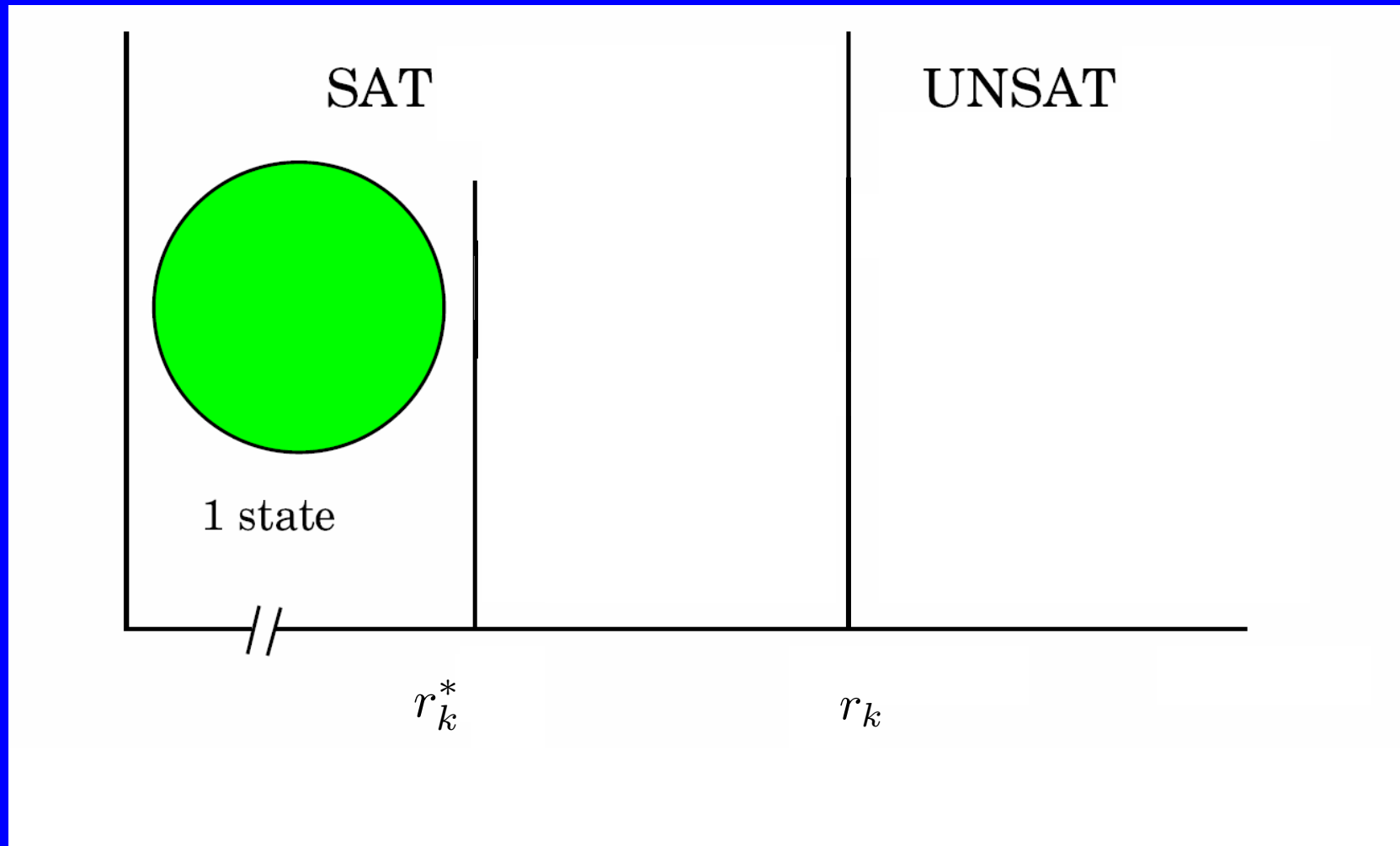
# Statistical Physics



SAT          UNSAT

$r_k$

# Statistical Physics

# Statistical Physics

# Statistical Physics

# Statistical Physics

# Sampling satisfying assignments

(thought experiment)

- Approximate the fraction $p_i$ of satisfying truth assignments in which variable $x_i$ takes value 1.

- Set $x_i$ to 1 with probability $p_i$ and simplify.

# Sampling satisfying assignments

(thought experiment)

- Approximate the fraction $p_i$ of satisfying truth assignments in which variable $x_i$ takes value 1.

- Set $x_i$ to 1 with probability $p_i$ and simplify.

**Locally:**

# Sampling satisfying assignments

(thought experiment)

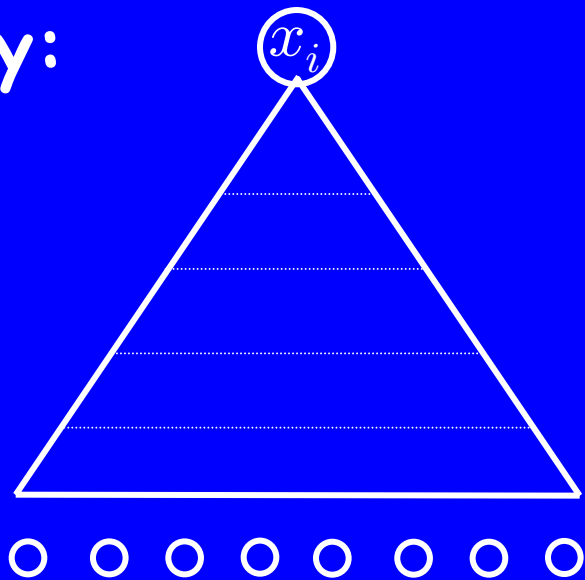- **Approximate** the fraction $p_i$ of satisfying truth assignments in which variable $x_i$ takes value 1.

- Set $x_i$ to 1 with probability $p_i$ and simplify.

**Locally:**



**Given boundary $\Lambda$:** compute $\mathrm{p}_\Lambda$

$$p_i = \sum_\Lambda p_\Lambda \times \mathrm{Ext}(\Lambda)$$

# Hope

- The variables in the boundary of the tree are "far apart in the graph" (if we remove the tree).
- Therefore, they should be uncorrelated; in which case "we can compute".

e.g., LDPC codes

# Hope

- The variables in the boundary of the tree are "far apart in the graph" (if we remove the tree).
- Therefore, they should be uncorrelated; in which case "we can compute".

e.g., LDPC codes

## But if clustering exists...

- The marginals are NOT uncorrelated.

- Clusters with many frozen variables induce "long-range" correlations.

# Rigorizing the 1-RSB picture

We **prove** that at $t_k \sim \dfrac{2^k}{k} \log k$

- Exponentially many clusters appear
- They are far apart from one another
- They have small diameter
- Many variables are frozen in each

# Rigorizing the 1-RSB picture

We **prove** that at $t_k \sim \dfrac{2^k}{k} \log k$

- Exponentially many clusters appear
- They are far apart from one another
- They have small diameter
- Many variables are frozen in each

---

Contrast: set of solutions is "convex" up to

$$\sim \frac{2^k}{k}$$

# Definitions

For any formula $F$:

-Let $\mathcal{S}(F)$ be the set of satisfying assignments of $F$.

-Let $C_1, C_2, \ldots$ be the connected components (clusters)

of $\mathcal{S}(F)$.   (Adjacent = Hamming distance 1)

-Let the label of $C$ be its projection $\ell(C) \in \{0, 1, *\}^n$.

-If $\ell_i(C) \in \{0, 1\}$ we say that $x_i$ is frozen in $C$.

Two quick observations:

- Labels are "lossless" for cubes.

- The label of $C$ can be "all-stars" already with $|C|=n$.

# A majority of frozen variables

**Theorem.** *For every $k \geq 9$ and*

$$r > c_k = \frac{4}{5} \, 2^k \ln 2 \, (1 + o(1)),$$

*w.h.p. in every cluster the majority of variables are frozen.*

# Nearly everything freezes

**Theorem.** *For every $\epsilon > 0$ and all $k \geq k_0(\epsilon)$, there exists $c_k^\epsilon < r_k$, such that w.h.p. in every cluster at least $(1 - \epsilon) \cdot n$ variables are frozen.*

# Nearly everything freezes

**Theorem.** *For every $\epsilon > 0$ and all $k \geq k_0(\epsilon)$, there exists $c_k^\epsilon < r_k$, such that w.h.p. in every cluster at least $(1 - \epsilon) \cdot n$ variables are frozen.*

*As $k$ grows,*

$$\frac{c_k^\epsilon}{2^k \ln 2} \rightarrow \frac{1}{1 + \epsilon(1 - \epsilon)}$$

# Coarsening

**Definition.** *A variable $x_i$ is* **free** *in $x \in \{0, 1, *\}^n$ if in every clause containing $x_i, \overline{x}_i$ there is some other satisfied literal or $*$.*

# Coarsening

**Definition.** *A variable $x_i$ is **free** in $x \in \{0, 1, *\}^n$ if in every clause containing $x_i, \overline{x}_i$ there is some other satisfied literal or $*$.*

Repeat until fixed point: set all free variables to $*$ .

# Coarsening

**Definition.** *A variable $x_i$ is **free** in $x \in \{0, 1, *\}^n$ if in every clause containing $x_i, \overline{x}_i$ there is some other satisfied literal or $*$.*

Repeat until fixed point: set all free variables to $*$ .

1. All $\sigma$ in $C$ have the same fixed point, called cover($C$).
2. label($C$) $\preceq$ cover($C$) deterministically.

# Proof

- Let X be the number of satisfying assignments whose cover (fixed point) is "all-∗".  (Call them "coreless".)

# Proof

- Let X be the number of satisfying assignments whose cover (fixed point) is "all-∗". (Call them "coreless".)

$$\mathbf{E}[X] = \sum_{\sigma} \Pr[\sigma \ is \ coreless \mid \sigma \ is \ satisfying] \times \Pr[\sigma \ is \ satisfying]$$

$$= 2^n \cdot \left(1 - \frac{1}{2^k}\right)^{rn} \cdot \Pr[\mathbf{0} \ is \ coreless \mid \mathbf{0} \ is \ satisfying]$$

# Proof

● Let X be the number of satisfying assignments whose cover (fixed point) is "all-$*$".  (Call them "coreless".)

$$\mathbf{E}[X] = \sum_{\sigma} \Pr[\sigma \ is \ coreless \mid \sigma \ is \ satisfying] \times \Pr[\sigma \ is \ satisfying]$$

$$= 2^n \cdot \left(1 - \frac{1}{2^k}\right)^{rn} \cdot \Pr[\mathbf{0} \ is \ coreless \mid \mathbf{0} \ is \ satisfying]$$

● Conditioning on "0 *is satisfying*" is easy
● Relevant clauses = uniquely-satisfied clauses
● Similar to hypergraph core computation

# Proof

- Let X be the number of satisfying assignments whose cover (fixed point) is "all-*".  (Call them "coreless".)

$$
\begin{aligned}
\mathbf{E}[X] &= \sum_{\sigma} \Pr[\sigma \text{ is coreless} \mid \sigma \text{ is satisfying}] \times \Pr[\sigma \text{ is satisfying}] \\[1em]
&= 2^n \cdot \left(1 - \frac{1}{2^k}\right)^{rn} \cdot \Pr[\mathbf{0} \text{ is coreless} \mid \mathbf{0} \text{ is satisfying}] \\[1em]
&< \left[2 \cdot \left(1 - \frac{1}{2^k}\right)^r \cdot e^{-f(r)}\right]^n
\end{aligned}
$$

# Proof

- Let X be the number of satisfying assignments whose cover (fixed point) is "all-*".  (Call them "coreless".)

$$\mathbf{E}[X] \;=\; \sum_\sigma \Pr[\sigma \ is \ coreless \mid \sigma \ is \ satisfying] \times \Pr[\sigma \ is \ satisfying]$$

$$=\; 2^n \cdot \left(1 - \frac{1}{2^k}\right)^{rn} \cdot \Pr[\mathbf{0} \ is \ coreless \mid \mathbf{0} \ is \ satisfying]$$

$$<\; \left[2 \cdot \left(1 - \frac{1}{2^k}\right)^r \cdot e^{-f(r)}\right]^n$$

$$\Pr[\mathbf{0} \ is \ coreless \mid \mathbf{0} \ is \ satisfying] \;=\; \begin{cases} 1 - o(1) & \text{if } r < t_k \\[2mm] o(1) & \text{if } r > t_k \end{cases}$$

# Proof

- Let X be the number of satisfying assignments whose cover (fixed point) is "all-$*$". (Call them "coreless".)

$$\mathbf{E}[X] = \sum_\sigma \Pr[\sigma \text{ is coreless} \mid \sigma \text{ is satisfying}] \times \Pr[\sigma \text{ is satisfying}]$$

$$= 2^n \cdot \left(1 - \frac{1}{2^k}\right)^{rn} \cdot \Pr[\mathbf{0} \text{ is coreless} \mid \mathbf{0} \text{ is satisfying}]$$

$$< \left[2 \cdot \left(1 - \frac{1}{2^k}\right)^r \cdot e^{-f(r)}\right]^n$$

$$t_k \sim \frac{2^k}{k} \log k$$

# Contiguity of Planted & Random

We will create two formulas with n variables and m=rn clauses, where $r < 2^k \ln 2 - k$:

# Contiguity of Planted & Random

We will create two formulas with n variables and m=rn clauses, where $r < 2^k \ln2 - k$:

- F is generated by selecting the m clauses at random, among all possible clauses.

# Contiguity of Planted & Random

We will create two formulas with n variables and m=rn clauses, where $r < 2^k \ln2 - k$:

- F is generated by selecting the m clauses at random, among all possible clauses.
- G is generated by:
    - Selecting a random τ in $\{0,1\}^n$.
    - Selecting m clauses compatible with τ at random.

# Contiguity of Planted & Random

We will create two formulas with n variables and m=rn clauses, where $r < 2^k \ln 2 - k$:

- F is generated by selecting the m clauses at random, among all possible clauses.
- G is generated by:
  - Selecting a random $\tau$ in $\{0,1\}^n$.
  - Selecting m clauses compatible with $\tau$ at random.

Let $\sigma$ be a random satisfying assignment of F (if one exists). The pairs $(\sigma,F)$ and $(\tau,G)$ are statistically indistinguishable.

# Summary

- Much before disappearing solutions form clusters:
  - Relatively small
  - Far apart
  - Exponentially many
- "Error-correcting-code with fuzz"

# Summary

- Much before disappearing solutions form clusters:
  - Relatively small
  - Far apart
  - Exponentially many
- "Error-correcting-code with fuzz"

---

- Frozen variables -> long range correlations -> cause naive local algorithms to fail.

# Summary

- Much before disappearing solutions form clusters:
    - Relatively small
    - Far apart
    - Exponentially many
- "Error-correcting-code with fuzz"

---

- Frozen variables -> long range correlations -> cause naive local algorithms to fail.

> Influence propagation without gadgets.