# CIS 700:
# "algorithms for Big Data"

# Lecture 1: Intro

Slides at http://grigory.us/big-data-class.html
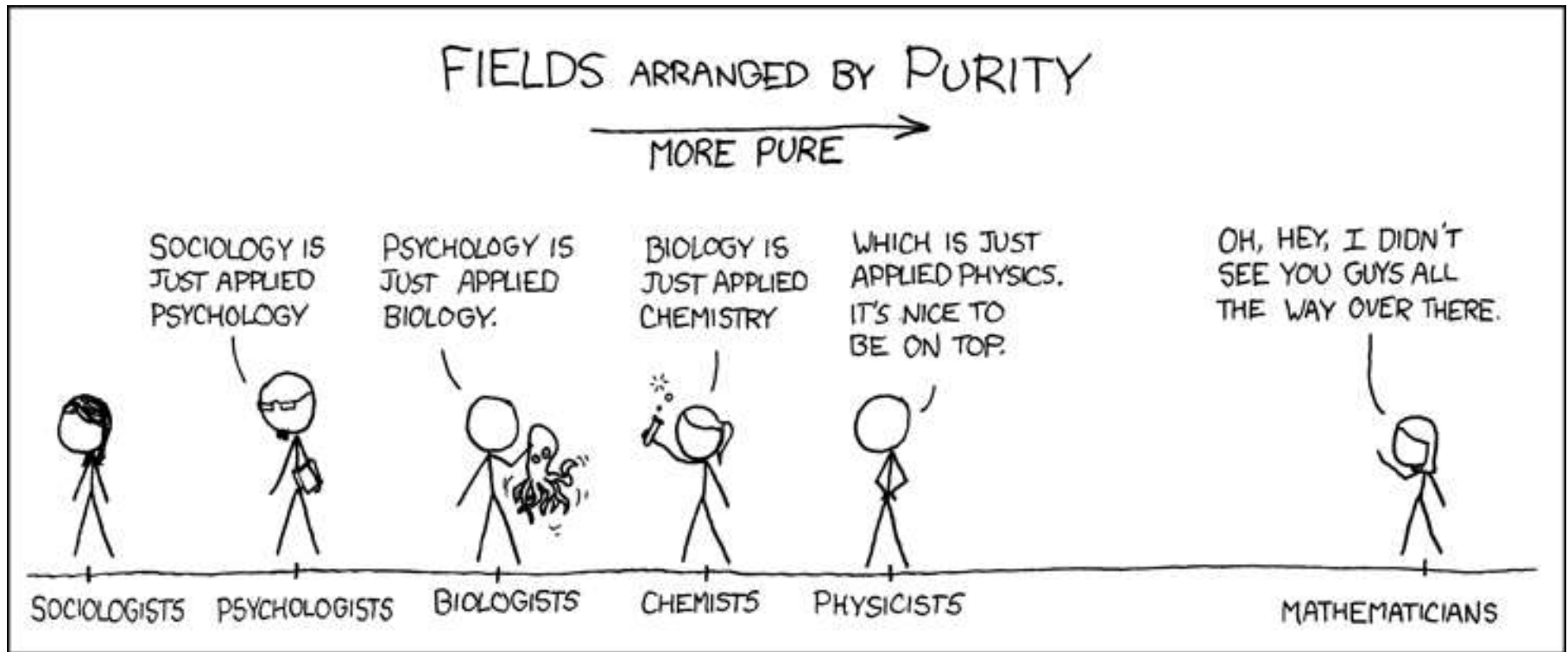
## Grigory Yaroslavtsev
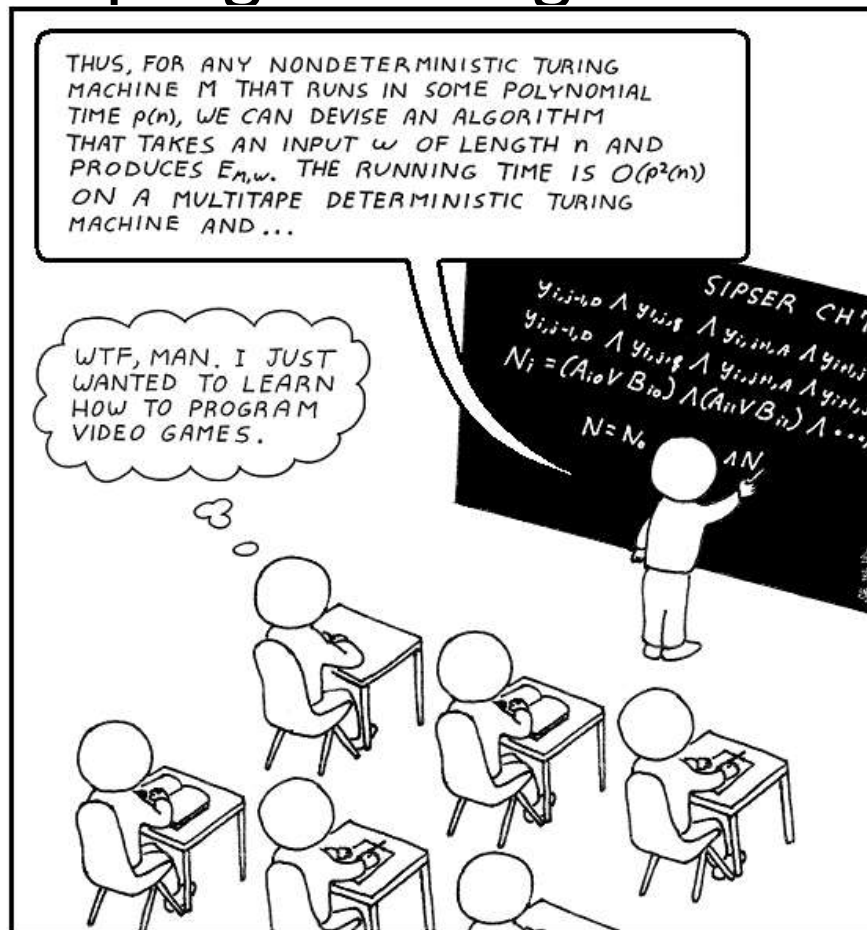
## http://grigory.us

# Disclaimers

- A lot of Math!

# Disclaimers

- (Almost) no programming!

# Class info

- MW 10:30 – 12:00, Towne 307
- Grading:
  - 1-2 homework assignments (40%)
  - Project (60%)
- Office hours by appointment
- Slides will be posted

# What is this class about?

- Not about the band (https://en.wikipedia.org/wiki/Big_Data_(band))

# What is this class about?

- The four V's: **volume**, **velocity**, variety, veracity

- **Volume:** "Big Data" = too big to fit in RAM
  – Today 16GB $\approx 100\$$ => "big" starts at terabytes

- **Velocity:** real-time
  – Doesn't fit in RAM + has to be processed on the fly

- **N** = size of data, time and memory o(**N**)

- o(**N**): $(1)$, $(\log N)$, $(\quad)$ where $0 <\quad< 1$

# Getting hands dirty

- Cloud computing platforms (all offer free trials):
    - Amazon EC2 (1 CPU/12mo)
    - Microsoft Azure ($200/1mo)
    - Google Compute Engine ($200/2mo)

- Distributed Google Code Jam
    - First time in 2015:

        https://code.google.com/codejam/distributed_index.html

    - Caveats:
        - Very basic aspects of distributed algorithms (few rounds)
        - Small data ($\sim 1$      , with hundreds MB RAM)
        - Fast query access ($\sim 0.01$      per request), "data with queries"

# Outline

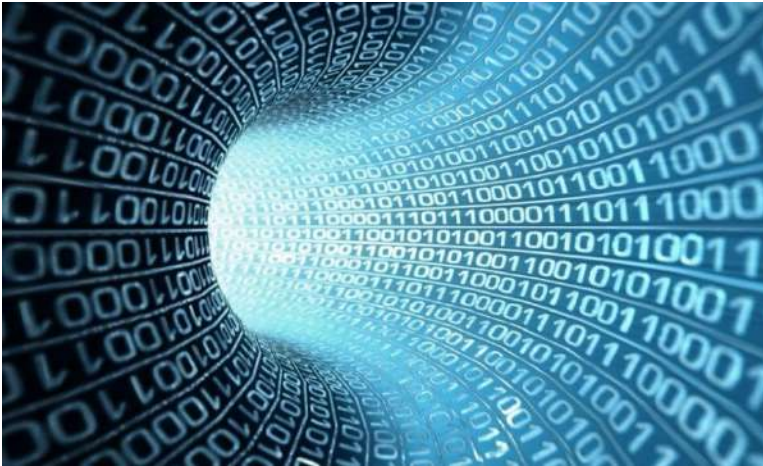- Part 1: Streaming Algorithms

Highlights:

- Approximate counting
- # Distinct Elements, Hyperloglog
- Median
- Frequency moments
- Heavy hitters
- Graph sketching

# Outline



- Part 2: Algorithms for numerical linear algebra

    Highlights:
    - Dimension reduction
    - Nearest neighbor search
    - Linear sketching
    - Linear regression
    - Low rank approximation

# Outline

- Part 3: Massively Parallel Algorithms

  Highlights:
  - Computational Model
  - Sorting (Terasort)
  - Connectivity, MST
  - Filtering dense graphs
  - Euclidean MST

# Outline

- Part 4: Sublinear Time Algorithms

Highlights:
- "Data with queries"
- Sublinear approximation
- Property Testing
- Testing images, sortedness, connectedness
- Testing noisy data

# Today

# Puzzles

You see a sequence of values $a_1, ..., a_m$, arriving one by one:

- (**Easy, "Find a missing player"**)
  - If all $a_i$ are different and have values between $1$ and $m + 1$, which value is missing?
  - You have $O(\log m)$ space

- Example:
  - There are 11 soccer players with numbers 1, ..., 11.
  - You see 10 of them one by one, which one is missing? You can only remember a single number.

1

8

5

11

3

9

2

6

7

4

# Which number was missing?

# Puzzle #1

You see a sequence of values $x_1, ..., x_n$, arriving one by one:

- (**Easy, "Find a missing player"**)
  - If all $x_i$ are different and have values between $1$ and $n + 1$, which value is missing?
  - You have $O(\log n)$ space
- Example:
  - There are 11 soccer players with numbers 1, ..., 11.
  - You see 10 of them one by one, which one is missing? You can only remember a single number.

# Puzzle #2

You see a sequence of values $_1, ..., $ , arriving one by one:

- (**Harder, "Keep a random team"**)
  - How can you maintain a uniformly random sample of values out of those you have seen so far?
  - You can store exactly items at any time

- Example:
  - You want to have a team of 11 players randomly chosen from the set you have seen.
  - Players arrive one at a time and you have to decide whether to keep them or not.

# Puzzle #3

You see a sequence of values $v_1, ..., v_n$, arriving one by one:

- (**Very hard, "Count the number of players"**)
  - What is the total number of values up to error $\pm \epsilon$?
  - You have $O(\log\log n / \epsilon^2)$ space and can be completely wrong with some small probability

# Puzzles

You see a sequence of values $x_1, \ldots, x_n$, arriving one by one:

- (**Easy, "Find a missing player"**)
  - If all $x_i$'s are different and have values between $1$ and $n + 1$, which value is missing?
  - You have $O(\log n)$ space
- (**Harder, "Keep a random team"**)
  - How can you maintain a uniformly random sample of $k$ values out of those you have seen so far?
  - You can store exactly $k$ items at any time
- (**Very hard, "Count the number of players"**)
  - What is the total number of values up to error $\pm \varepsilon$?

# Part 1: Probability 101

"The bigger the data the better you should know your Probability"

- Basic Probability:
  - Probability, events, random variables
  - Expectation, variance / standard deviation
  - Conditional probability, independence, pairwise independence, mutual independence

# Expectation

- = random variable with values $_1, ..., \quad , ...$

- Expectation $\mathbb{E}[\quad]$

$$\mathbb{E}[\quad] = \sum_{=1}^{\infty} x_i \cdot \Pr[\quad = \quad ]$$

- Properties (linearity):

$$\mathbb{E}[\quad] = \mathbb{E}[\quad]$$

$$\mathbb{E}[\quad + \quad] = \mathbb{E}[\quad] + \mathbb{E}[\quad]$$

- Useful fact: if all $\geq 0$ and integer then

$$\sum^{\infty}$$

# Variance

- Variance $\quad [ \quad ] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])^2]$

$$[ \quad ] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])^2] =$$
$$= \mathbb{E}\left[ \quad ^2 - 2 \quad \cdot \mathbb{E}[\mathbf{X}] + \mathbb{E}[\mathbf{X}]^2\right]$$
$$= \mathbb{E}[ \quad ^2] - \textcolor{blue}{2\mathbb{E}[ \quad \cdot \mathbb{E}[\mathbf{X}]]} + \textcolor{green}{\mathbb{E}[\mathbb{E}[\mathbf{X}]^2]}$$

- $\mathbb{E}[X]$ is some fixed value (a constant)
- $\textcolor{blue}{2\, \mathbb{E}[ \quad \cdot \mathbb{E}[\mathbf{X}]] = 2\, \mathbb{E}[\mathbf{X}] \cdot \mathbb{E}[\mathbf{X}] = 2\, \mathbb{E}^2[ \quad ]}$
- $\textcolor{green}{\mathbb{E}[\mathbb{E}[\mathbf{X}]^2] = \mathbb{E}^2[\mathbf{X}]}$

# Independence

- Two random variables    and    are **independent** if and only if (iff) for every   ,   :

$$\Pr\left[\quad = \quad , \quad = \quad\right] = \Pr\left[\quad = \quad\right] \cdot \Pr\left[\quad = \quad\right]$$

- Variables   $_1, \ldots,$   are **mutually independent** iff

$$\Pr\left[\quad = \quad _1, \ldots, \quad = \quad\right] = \prod_{=1} \Pr\left[\quad = \quad\right]$$

- Variables   $_1, \ldots,$   are **pairwise independent** iff for all pairs i,j

$$\Pr\left[\quad = \quad , \quad = \quad\right] = \Pr\left[\quad = \quad\right] \Pr\left[\quad = \quad\right]$$

# Conditional Probabilities

- For two events $E_1$ and $E_2$:

$$\Pr\left[E_2 \mid E_1\right] = \frac{\Pr\left[E_1 \cap E_2\right]}{\Pr\left[E_1\right]}$$

- If two random variables (r.vs) are independent

$$\Pr\left[X_2 = x_2 \mid X_1 = x_1\right]$$

$$= \frac{\Pr\left[X_1 = x_1 \cap X_2 = x_2\right]}{\Pr\left[X_1 = x_1\right]} \text{ (by definition)}$$

$$= \frac{\Pr\left[X_1 = x_1\right]\Pr\left[X_2 = x_2\right]}{\Pr\left[X_1 = x_1\right]} \text{ (by independence)}$$

# Union Bound

For any events $A_1, ..., A_k$:

$$\Pr[A_1 \cup A_2 \cup ... \cup A_k] \leq \Pr[A_1] + \Pr[A_2] + ... + \Pr[A_k]$$

- **Pro**: Works even for dependent variables!
- **Con**: Sometimes very loose, especially for **mutually**

# Independence and Linearity of Expectation/Variance

- Linearity of expectation (even for dependent variables!):

$$\mathbb{E}\left[\sum_{=1}\right] = \sum_{=1}\mathbb{E}[\quad]$$

- Linearity of variance (only for **pairwise independent** variables!)

$$\left[\sum_{=1}\right] = \sum_{=1}\quad[\quad]$$

# Part 2: Inequalities

- Markov inequality
- Chebyshev inequality
- Chernoff bound

# Markov's Inequality

- For every $> 0:$ $Pr\left[\ \geq\ \left[\ \right]\right] \leq \dfrac{1}{}$

- **Proof (by contradiction)** $Pr\left[\ \geq\ \left[\ \right]\right] \color{red}{>} \dfrac{1}{}$

$$\left[\ \right] = \sum\ \cdot \Pr\left[\ \right] = \ ]\qquad \text{(by definition)}$$

$$\geq \sum_{=\ \left[\ \right]}^{\infty}\ \cdot \Pr\left[\ =\ \right]\qquad \text{(pick only some i's)}$$

# Markov's Inequality

- For every $> 0$: $Pr[ \geq [ ]] \leq \dfrac{1}{}$

- **Corollary** $(c' = [ ])$ :

For every $' > 0$: $Pr[ \geq '] \leq \dfrac{[ ]}{'}$

- **Pro**: always works!

- **Cons**:
  - Not very precise
  - Doesn't work for the lower tail: $Pr[ \leq [ ]]$

# Chebyshev's Inequality

- For every $> 0$:

$$\Pr\left[\,|\quad - \quad [\quad]\,| \geq \sqrt{\quad[\quad]}\,\right] \leq \frac{1}{2}$$

- Proof:

$$\Pr\left[\,|\quad - \quad[\quad]\,| \geq \sqrt{\quad[\quad]}\,\right]$$

$$= \Pr\left[\,|\quad - \quad[\quad]|^2 \geq \,^2\quad[\quad]\,\right] \qquad \text{(by squaring)}$$

$$= \Pr\left[\,|\quad - \quad[\quad]|^2 \geq \,^2\,[\,|\quad - \quad[\quad]|^2\,]\,\right] \text{ (def. of Var}$$

$$\leq \frac{1}{2} \qquad \text{(by Markov's inequality)}$$

# Chebyshev's Inequality

- For every $\quad > 0$:

$$\Pr\left[\lvert \quad - \quad[\quad] \rvert \geq \sqrt{\quad[\quad]}\right] \leq \frac{1}{\quad^2}$$

- **Corollary** ($\quad' = \sqrt{\quad[\quad]}$):

For every $\quad' > 0$:

$$\Pr\left[\lvert \quad - \quad[\quad] \rvert \geq \quad'\right] \leq \frac{\quad[\quad]}{\quad'^2}$$

# Chernoff bound

- Let $X_1 \ldots X_n$ be independent and identically distributed r.vs with range [0,1] and expectation $\mu$.

- Then if $X = \frac{1}{n}\sum X_i$ and $1 > \varepsilon > 0$,

$$\Pr\left[\,|X - \mu| \geq \varepsilon\,\right] \leq 2 \exp\left(-\frac{n\varepsilon^2}{3}\right)$$

# Chernoff bound (corollary)

- Let $X_1 \dots X_n$ be independent and identically distributed r.vs with range [0, **c**] and expectation $\mu$.

- Then if $X = \dfrac{1}{n} \displaystyle\sum X_i$ and $1 > \delta > 0,$

$$\Pr\left[\,|X - \mu| \geq \delta\mu \,\right] \leq 2 \exp\left( -\frac{n\mu\delta^2}{3c} \right)$$

# Chernoff v.s Chebyshev

Large values of t is exactly what we need!

Let $X_1 \ldots$ be independent and identically distributed r.vs with range [0,1] and expectation $\mu$. Let $Z = \frac{1}{n}\sum X_i$.

- Chebyshev: $\Pr\left[|Z - \mu| \geq t\right] = O\left(\frac{1}{n}\right)$

- Chernoff: $\Pr\left[|Z - \mu| \geq t\right] = e^{-\Omega(n)}$

So is Chernoff always better for us?
- Yes, if we have i.i.d. variables

# Answers to the puzzles

You see a sequence of values $x_1, ..., x_n$, arriving one by one:

- (**Easy**)
  - If all $x_i$ are different and have values between $1$ and $n + 1$, which value is missing?
  - You have $O(\log n)$ space

  - **Answer**: missing value $= \sum_{i=1}^{n+1} i - \sum_{i=1}^{n} x_i$

- (**Harder**)
  - How can you maintain a uniformly random sample of values out of those you have seen so far?

# Part 3: Morris's Algorithm

- (**Very hard, "Count the number of players"**)
  - What is the total number of values up to error $\pm \quad ?$
  - You have $( \log \log \quad / \quad ^2)$ space and can be completely wrong with some small probability

# Morris's Algorithm: Alpha-version

Maintains a counter    using $log\;log\;f0$ bits

- Initialize    to 0
- When an item arrives, increase X by 1 with probability $\dfrac{1}{2}$
- When the stream is over, output $2 - 1$

Claim:    $\left[2\;\right] = + 1$

# Morris's Algorithm: Alpha-version

Maintains a counter using $\log \log f_0$ bits

- Initialize to 0, when an item arrives, increase X by 1 with probability $\dfrac{1}{2}$

Claim: $[2^{\phantom{X}}] = \phantom{X} + 1$

- Let the value after seeing items be

$$[2^{\phantom{X}}] = \sum_{=0}^{\infty} \Pr[_{\phantom{f_0}} {}_{-1} = \phantom{X}] \; [2^{\phantom{X}} \mid \phantom{X}_{-1} = \phantom{X}]$$

$$\sum^{\infty} \phantom{\Pr} \left( \frac{1}{\phantom{2}} \phantom{2} + 1 \left( \phantom{X} \frac{1}{\phantom{X}} \right) \phantom{2} \right)$$

# Morris's Algorithm: Alpha-version

Maintains a counter using $\log\log f_0$ bits

- Initialize to 0, when an item arrives, increase X by 1 with probability $\dfrac{1}{2}$

Claim: $\left[2^{2}\right] = \dfrac{3}{2}\,{}^{2} + \dfrac{3}{2}\, + 1$

$$\left[2^{2}\right] = \sum_{=0}^{\infty} \Pr[2 \quad {}_{-1} = \ ] \left[2^{2} \,/2 \quad {}_{-1} = \ \right]$$

$$= \sum^{\infty} \Pr[2 \quad {}_{-1} = \ ]\left(\frac{1}{-}4\ {}^{2} + \left(1 - \frac{1}{-}\right){}^{2}\right)$$

# Morris's Algorithm: Alpha-version

Maintains a counter using $\log\log f_0$ bits

- Initialize to 0, when an item arrives, increase X by 1 with probability $\frac{1}{2}$

- $[2^{\ }] = n + 1, \qquad [2^{\ }] = (\ ^2)$

- Is this good?

# Morris's Algorithm: Beta-version

Maintains counters $c_1, ..., c_k$ using $\log\log f_0$ bits for each

- Initialize $c_i'$ to 0, when an item arrives, increase each $c_i$ by 1 independently with probability $\dfrac{1}{2^{c_i}}$

- Output $Z = \dfrac{1}{k}\left(\displaystyle\sum_{i=1}^{k} 2^{c_i} - 1\right)$

- $E[2^{c_i}] = n + 1, \quad E[2^{c_i}] = \left(\begin{array}{c} 2 \end{array}\right)$

$\left(\dfrac{1}{k}\sum\right) \quad \left(\begin{array}{c} 2 \end{array}\right)$

# Morris's Algorithm: Beta-version

Maintains counters $c_1, ..., c_q$ using $\log\log f_0$ bits for each

- Output $Z = \frac{1}{q}(\sum_{i=1}^{q} 2^{c_i} - 1)$

- $E[Z] = E\left(\frac{1}{q}\sum_{i=1}^{q} 2^{c_i} - 1\right) = E\left(\frac{2^{c}}{}\right)$

- Claim: If $q \geq \frac{1}{\varepsilon^2}$ then $\Pr\left[|Z - f_0| > \varepsilon f_0\right] < 1/3$

# Morris's Algorithm: Final

- What if I want the probability of error to be really small, i.e. $\Pr\left[\,|\quad - \quad| > \quad\right] \le \;?$

- Same Chebyshev-based analysis: $\quad = \quad\left(\dfrac{1}{2}\right)$

- Do these steps $\quad = \quad\left(\log\dfrac{1}{-}\right)$ times

  independently in parallel and output the median answer.

  $$\left(\log\log\quad \cdot \log\dfrac{1}{-}\right)$$

- Total space:

# Morris's Algorithm: Final

- Do these steps $= \left( \log \frac{1}{\,} \right)$ times independently in parallel and output the median answer .

Maintains counters $^1, \ldots,$ using $\log \log f_0$ bits for each

- Initialize $'$ to 0, when an item arrives, increase each by 1 independently with probability $\frac{1}{2}$

# Morris's Algorithm: Final Analysis

Claim: $\Pr\left[\left|\quad - \quad\right| > \quad\right] \leq$

- Let $\quad$ be an indicator r.v. for the event that $\left|\quad - \quad\right| \leq \quad$, where $\quad$ is the i-th trial.

- Let $\quad = \sum\quad$.

- $\Pr\left[\left|\quad - \quad\right| > \quad\right] \leq \Pr\left[\quad \leq \dfrac{\quad}{2}\right]$

# Thank you!

- Questions?
- **Next time**:
  - More streaming algorithms