

CS-GY 6763: Lecture 3

Exponential Concentration Inequalities, Fingerprinting

NYU Tandon School of Engineering, Prof. Christopher Musco

Lemma (Chebyshev's Inequality)

Let X be a random variable with expectation $\mathbb{E}[X]$ and variance $\sigma^2 = \text{Var}[X]$. Then for any $k > 0$,

$$\Pr[|X - \mathbb{E}[X]| \geq k \cdot \sigma] \leq \frac{1}{k^2}$$

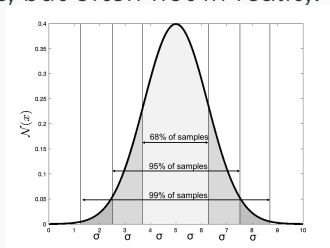
One application: Proved that if you throw n balls into n bins, the maximum loaded bin has $O(\sqrt{n})$ balls. We used Chebyshev's + _____.



This lecture, we'll prove a bound of $O(\log n)$ using stronger tools.

Motivating question: Is Chebyshev's Inequality tight?

It is the worst case, but often not in reality.



68-95-99 rule for Gaussian bell-curve. $X \sim N(0, \sigma^2)$

Chebyshev's Inequality:

$$\Pr(|X - \mathbb{E}[X]| \geq 1\sigma) \leq 100\%$$

$$\Pr(|X - \mathbb{E}[X]| \geq 2\sigma) \leq 25\%$$

$$\Pr(|X - \mathbb{E}[X]| \geq 3\sigma) \leq 11\%$$

$$\Pr(|X - \mathbb{E}[X]| \geq 4\sigma) \leq 6\%.$$

Truth:

$$\Pr(|X - \mathbb{E}[X]| \geq 1\sigma) \approx 32\%$$

$$\Pr(|X - \mathbb{E}[X]| \geq 2\sigma) \approx 5\%$$

$$\Pr(|X - \mathbb{E}[X]| \geq 3\sigma) \approx 1\%$$

$$\Pr(|X - \mathbb{E}[X]| \geq 4\sigma) \approx .01\%$$

GAUSSIAN CONCENTRATION

For $X \sim \mathcal{N}(\mu, \sigma^2)$:

$$\Pr[X = \mu \pm x] \sim \frac{1}{\sigma\sqrt{2\pi}} e^{-x^2/2\sigma^2}$$

Lemma (Gaussian Tail Bound)

For $X \sim \mathcal{N}(\mu, \sigma^2)$:

$$\Pr[|X - \mathbb{E}X| \geq k \cdot \sigma] \leq 2e^{-k^2/2}.$$

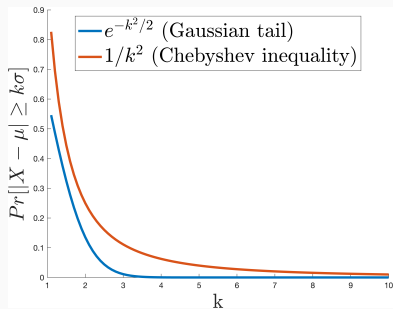
Compare this to:

Lemma (Chebyshev's Inequality)

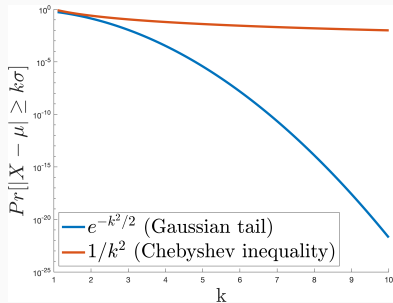
For $X \sim \mathcal{N}(\mu, \sigma^2)$:

$$\Pr[|X - \mathbb{E}X| \geq k \cdot \sigma] \leq \frac{1}{k^2}$$

GAUSSIAN CONCENTRATION



Standard y-scale.



Logarithmic y-scale.

Takeaway: Gaussian random variables concentrate much tighter around their expectation than variance alone (i.e. Chebyshev's inequality) predicts.

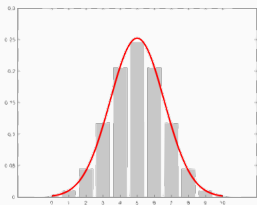
Why does this matter for algorithm design?

CENTRAL LIMIT THEOREM

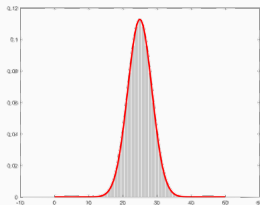
Theorem (CLT – Informal)

Any sum of *mutually independent, (identically distributed)* r.v.'s X_1, \dots, X_k with mean μ and finite variance σ^2 converges to a Gaussian r.v. with mean $k \cdot \mu$ and variance $k \cdot \sigma^2$, as $k \rightarrow \infty$.

$$S = \sum_{i=1}^n X_i \implies \mathcal{N}(k \cdot \mu, k \cdot \sigma^2).$$



(a) Distribution of # of heads after 10 coin flips, compared to a Gaussian.



(b) Distribution of # of heads after 50 coin flips, compared to a Gaussian.

Recall:

Definition (Mutual Independence)

Random variables X_1, \dots, X_k are mutually independent if, for all possible values v_1, \dots, v_k ,

$$\Pr[X_1 = v_1, \dots, X_k = v_k] = \Pr[X_1 = v_1] \cdot \dots \cdot \Pr[X_k = v_k]$$

Strictly stronger than pairwise independence.

EXERCISE

If I flip a fair coin 100 times, lower bound the chance I get between 30 and 70 heads?

For this problem, we will assume the limit of the CLT holds exactly – i.e., that this sum looks exactly like a Gaussian random variable.

Lemma (Gaussian Tail Bound)

For $X \sim \mathcal{N}(\mu, \sigma^2)$:

$$\Pr[|X - \mathbb{E}X| \geq k \cdot \sigma] \leq 2e^{-k^2/2}.$$

$2e^{-8} = .06\%$. Chebyshev's inequality gave a bound of 6.25%.

These back-of-the-envelope calculations can be made rigorous! Lots of different “versions” of bound which do so.

- Chernoff bound
- Bernstein bound
- Hoeffding bound
- ...

Different assumptions on random variables (e.g. binary vs. bounded), different forms (additive vs. multiplicative error), etc. **Wikipedia is your friend.**

Theorem (Chernoff Bound)

Let X_1, X_2, \dots, X_k be independent $\{0, 1\}$ -valued random variables and let $p_i = \mathbb{E}[X_i]$, where $0 < p_i < 1$. Then the sum $S = \sum_{i=1}^k X_i$, which has mean $\mu = \sum_{i=1}^k p_i$, satisfies

$$\Pr[S \geq (1 + \epsilon)\mu] \leq e^{\frac{-\epsilon^2 \mu}{2 + \epsilon}}.$$

and for $0 < \epsilon < 1$

$$\Pr[S \leq (1 - \epsilon)\mu] \leq e^{\frac{-\epsilon^2 \mu}{2}}.$$

Theorem (Chernoff Bound Corollary)

Let X_1, X_2, \dots, X_k be independent $\{0, 1\}$ -valued random variables and let $p_i = \mathbb{E}[X_i]$, where $0 < p_i < 1$. Let $S = \sum_{i=1}^k X_i$ and $\mathbb{E}[S] = \mu$. For $\epsilon \in (0, 1)$,

$$\Pr[|S - \mu| \geq \epsilon\mu] \leq 2e^{-\epsilon^2\mu/3}$$

Why does this look like the Gaussian tail bound of $\Pr[|S - \mu| \geq k \cdot \sigma] \lesssim 2e^{-k^2/2}$? What is $\sigma(S)$?

Theorem (Bernstein Inequality)

Let X_1, X_2, \dots, X_k be independent random variables with each $X_i \in [-1, 1]$. Let $\mu_i = \mathbb{E}[X_i]$ and $\sigma_i^2 = \text{Var}[X_i]$. Let $\mu = \sum_i \mu_i$ and $\sigma^2 = \sum_i \sigma_i^2$. Then, for $k \leq \frac{1}{2}\sigma$, $S = \sum_i X_i$ satisfies

$$\Pr[|S - \mu| > k \cdot \sigma] \leq 2e^{-k^2/4}.$$

Theorem (Hoeffding Inequality)

Let X_1, X_2, \dots, X_k be independent random variables with each $X_i \in [a_i, b_i]$. Let $\mu_i = \mathbb{E}[X_i]$ and $\mu = \sum_i \mu_i$. Then, for any $k > 0$, $S = \sum_i X_i$ satisfies:

$$\Pr[|S - \mu| > k] \leq 2e^{-\frac{k^2}{\sum_{i=1}^k (b_i - a_i)^2}}.$$

HOW ARE THESE BOUNDS PROVEN?

Variance is a natural measure of central tendency, but there are others.

q^{th} central moment: $\mathbb{E}[(X - \mathbb{E}X)^q]$

$q = 2$ gives the variance. Proof of Chebyshev's applies Markov's inequality to the random variable $(X - \mathbb{E}X)^2$.

Idea in brief: Apply Markov's inequality to $\mathbb{E}[(X - \mathbb{E}X)^q]$ for larger q , or more generally to $f(X - \mathbb{E}X)$ for some other non-negative function f . E.g., to $\exp(X - \mathbb{E}X)$.

EXERCISE

If I flip a fair coin 100 times, lower bound the chance I get between 30 and 70 heads?

Corollary of Chernoff bound: Let $S = \sum_{i=1}^k X_i$ and $\mu = \mathbb{E}[S]$. For $0 < \epsilon < 1$,

$$\Pr[|S - \mu| \geq \epsilon\mu] \leq 2e^{-\epsilon^2\mu/3}$$

Here $X_i = \mathbb{1}[i^{\text{th}} \text{ flip is heads}]$.

1.4%.

General Statement: Flip biased coin k times: i.e. the coin is heads with probability b . As long as $k \geq O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$,

$$\Pr[|\# \text{ heads} - b \cdot k| \geq \epsilon k] \leq \delta$$

Pay very little for higher probability – if you increase the number of coin flips by $4x$, δ goes from

$1/10 \rightarrow 1/100 \rightarrow 1/10000$

LOAD BALANCING

Recall: n jobs are distributed randomly to n servers using a hash function. Let S_i be the number of jobs sent to server i . What's the smallest B for which we can prove:

$$\Pr[\max_i S_i \geq B] \leq 1/10$$



Recall: Suffices to prove that, for any i , $\Pr[S_i \geq B] \leq 1/10n$:

$$\begin{aligned} \Pr[\max_i S_i \geq B] &= \Pr[S_1 \geq B \text{ or } \dots \text{ or } S_n \geq B] \\ &\leq \Pr[S_1 \geq B] + \dots + \Pr[S_n \geq B] \quad (\text{union bound}). \end{aligned}$$

Theorem (Chernoff Bound)

Let X_1, X_2, \dots, X_n be independent $\{0, 1\}$ -valued random variables and let $p_i = \mathbb{E}[X_i]$, where $0 < p_i < 1$. Then the sum $S = \sum_{j=1}^n X_j$, which has mean $\mu = \sum_{j=1}^n p_j$, satisfies

$$\Pr[X \geq (1 + \epsilon)\mu] \leq e^{\frac{-\epsilon^2 \mu}{2 + \epsilon}}.$$

Consider a single bin. Let $X_j = \mathbb{1}[\text{ball } j \text{ lands in that bin}]$.

$$\Pr[S \geq (1 + c \log n)\mu] \leq e^{\frac{-c^2 \log^2 n}{2 + c \log n}} \leq e^{\frac{-c \log^2 n}{2 \log n}} \leq e^{-.5c \log n} \leq \frac{1}{10n},$$

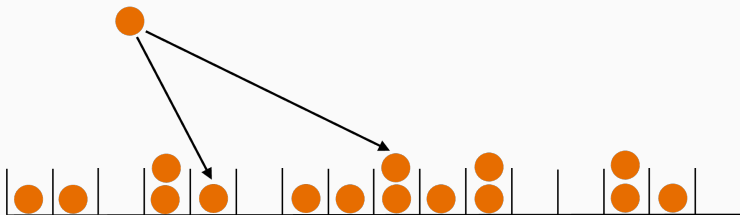
for sufficiently large c

So max load for randomized load balancing is $O(\log n)$! Best we could prove with Chebyshev's was $O(\sqrt{n})$.

POWER OF TWO CHOICES

Power of 2 Choices: Instead of assigning job to random server, choose 2 random servers and assign to the least loaded. With probability $1/10$ the maximum load is bounded by:

- (a) $O(\log n)$ (b) $O(\sqrt{\log n})$ (c) $O(\log \log n)$ (d) $O(1)$



BREAK

Recall from last class:

Definition (Universal hash function)

A random hash function $h : \mathcal{U} \rightarrow \{1, \dots, m\}$ is universal if, for any fixed $x, y \in \mathcal{U}$,

$$\Pr[h(x) = h(y)] \leq \frac{1}{m}.$$

Efficient construction: Let p be a prime number between $|\mathcal{U}|$ and $2|\mathcal{U}|$. Let a, b be random numbers in $0, \dots, p$, $a \neq 0$.

$$h(x) = [a \cdot x + b \pmod{p}] \pmod{m}$$

is universal.

We're not going to prove this, but this year I want to give a flavor for what tools are used.

One of the most famous applications of randomness in algorithm design.

Computational Problem: Given a number x , is x prime?

Recall:

- A number is prime if it can only be divided evenly by 1 and itself.
- The first few primes are 2, 3, 5, 7, 11, 13, 17, 19,
- Is 2023 prime?
- What about 49301977064557719291?

How would you design a generic algorithm to check?

Suppose we have an integer represented as a length n binary string.

$$x = 0110100010110101 \dots 1010001110$$

The naive prime checking algorithm runs in $O(2^n)$ time.

NYU Greene Super Computer has 2 petaFLOPS of throughput. When $n = 128$, would need 1 million Greene Super computers running for 1 million years to check if x is prime.

Miller-Rabin 1976, 1980: There is a randomized algorithm running in $O(n^3 \log_2(1/\delta))$ time that, with probability $1 - \delta$ determines if an n -bit integer x is prime.

- $n = 128$
- $\delta = 10^{-9}$ (chance of winning the Powerball Jackpot)
- $n^3 \log_2(1/\delta) \approx 60$ million operations.

Could check in $< .1$ second on my laptop.

This was a really big breakthrough!

Took over 20 more years to find a deterministic polynomial time primality test.

PRIMES is in P

Manindra Agrawal Neeraj Kayal
Nitin Saxena*

Department of Computer Science & Engineering
Indian Institute of Technology Kanpur
Kanpur-208016, INDIA
Email: {manindra,kayal,nitinsa}@iitk.ac.in

Abstract

We present an unconditional deterministic polynomial-time algorithm that determines whether an input number is prime or composite.

WHY DO PRIMES MATTER?

Basis for modern public-key cryptography.

Goal: Bob wants to send Alice an email. Wants to **encrypt** it in some way so that even if it is intercepted, no one can read it besides Alice.



```
fsxtwdrweodautbeiqxnhmwtnhzutzgf  
orzjfbryvokrtggbyeofchfyafekoyia  
gcdfndzqyunfwgwfivfvzvlrehcgxctj  
pcerkwdeysaruspsutecqvrbcetoiocg  
oqnajpsslcjdlcwbdcaxcglbtevfzhiu  
eutxubpmqnmxylonkaplmmqcvblbfltl  
dbjasqnsavwlgbazbmkbphezgeavtmet  
desrgtomtjrsxsrlshikycbpmhvincuhm  
ipzzwehmjnmriareccuxfqhjjezdsuqt  
zhnmyummpzimvydkvscsbtjrciquyfn
```

WHY DO PRIMES MATTER?

Basis for modern public-key cryptography.

Goal: Bob wants to send Alice an email. Wants to **encrypt** it in some way so that even if it is intercepted, no one can read it besides Alice.

Option 1: Share some sort of secret key/codebook in advance.



Impractical if you have a large number of uncoordinated senders and receivers.

Option 2: Create a 1-way lock box.



Anyone can deliver, only Alice can open/read the messages.

WHY DO PRIMES MATTER?

RSA cryptosystem (Rivest, Shamir, Adleman 1977):

- **Private key:** Two large (e.g. 128 bit) prime numbers p, q .
- **Public key:** Based on $z = p \times q$.

Even though checking if a number is prime can be done quickly, we do not have efficient algorithms for factoring numbers. E.g. for finding p, q based on z .¹

¹At least on classical computers we don't... different story on quantum computers.

How to find a 128 bit prime number p ? Use randomness, twice.

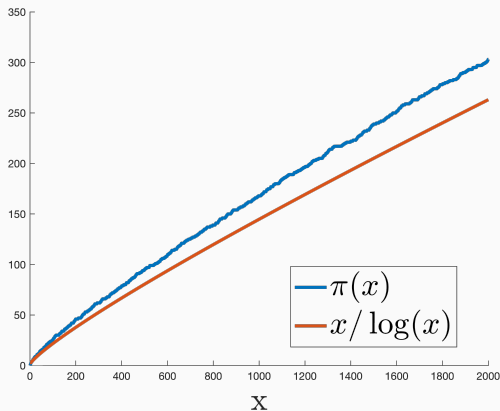
- Pick a random 128 bit number.
- Check if it's prime (using randomized primality test).
- If not, repeat.

Roughly how many tries do you expect this to take?

PRIME NUMBER THEOREM

Let $\pi(x)$ denote the number of primes less than some integer x . **Informally:**

$$\pi(x) \sim \frac{x}{\log(x)}$$



PRIME NUMBER THEOREM

Formally: For $x > 17$,

$$\frac{x}{\log(x)} \leq \pi(x) \leq \frac{x}{\log(x) - 4}$$

So if we select a random 128 bit number p , the chance that it is prime is great than:

$$\frac{1}{\log(2^{128})} \geq \frac{1}{90}$$

After a few hundred tries, we will almost definitely find a prime number. **In general, need $O(n)$ tries to find a prime with n bits.**

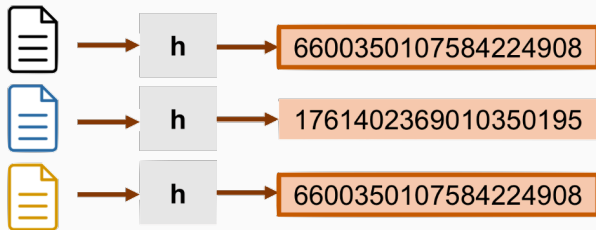
Finding large prime numbers is also a critical step in constructing efficiently computable universal hash.

Remainder of lecture: Discuss a simple but really important application of prime numbers to hashing.

FINGERPRINTING

Goal: Construct a compact “fingerprint” $h(f)$ for any given file f with two properties:

- The fingerprints $h(f_1)$ and $h(f_2)$ should be different with high probability if the contents of f_1 and f_2 differ at all.
- If the contents of f_1 and f_2 are identical, we should have $h(f_1) = h(f_2)$.



APPLICATIONS OF FINGER PRINTING

- Quickly check if two versions of the same file are identical (e.g. in version control systems like Git). Do not need to communicate the entire file between servers. Also used in webcaching and content delivery networks.
- Check that a file pieced together from multiple parts is not missing anything.

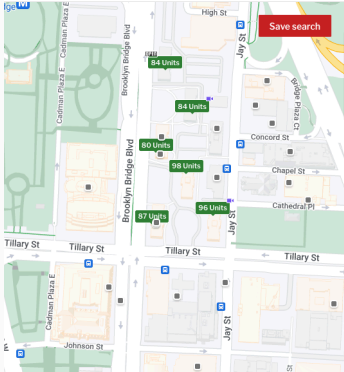
APPLICATIONS OF FINGER PRINTING



REDFIN Brooklyn... 1-844-759-7732 Buy • Rent • Sell • Redfin Premier Mortgage • Real Estate Agents • Feed Log In Sign Up

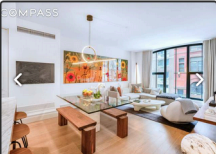
New York Homes for Sale [Market insights](#) | [City guide](#)



For sale ▾ Price ▾ Home type ▾ Beds / Baths ▾ All filters

17 homes • Sort: **Recommended** ▾ Photos Table

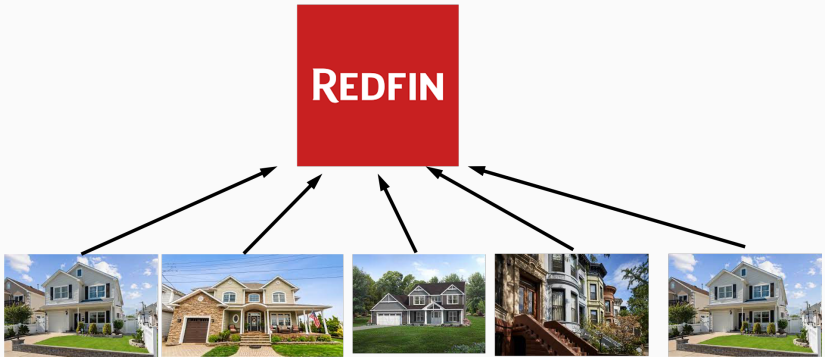


\$2,495,000  
3 Beds 2 Baths 1,429 Sq. Ft.
76 Schermerhorn St Unit 2-A, Brooklyn, NY 11201
Listing by Compass



\$825,000  
2 Beds 1 Bath 893 Sq. Ft.
225 ADAMS St Unit 15H, Brooklyn, NY 11201
Listing by Corcoran

APPLICATIONS OF FINGER PRINTING



Images from databases of local real estate agencies.

Fingerprints used as file names for the images to make sure we did not reupload new images that we already had, and to detect duplicate images and listings.

Goal: Construct a compact “fingerprint” function $h(f)$ such that:

- $h(f_1) \neq h(f_2)$ if $f_1 \neq f_2$ with high probability.

Ideally, length of $h(f_1)$ (i.e. the size of the integers hashed to) is much less than the file size.

Rabin Fingerprint (1981): Let file $f = 010 \dots 1101$ of length n be interpreted as an n bit integer. So something between 0 and 2^n .

Construct h randomly: Choose random prime number p between 2 and $tn \log(tn)$ for a constant t .

$$h(f) = f \pmod{p}.$$

How many bits does $h(f)$ take to store?

$$h(f) = f \pmod{p} \quad \text{for prime } p \in \{2, \dots, tn \log(tn)\}$$

Claim: If $f_1 \neq f_2$ then $h(f_1) = h(f_2)$ with probability $\leq \frac{2}{t}$.

Since our fingerprint only takes $O(\log n + \log t)$ space, we can set t to be super large, so effectively the probability of $h(f_1)$ and $h(f_2)$ colliding is negligible for all real-world applications.

E.g. set fingerprint length to $\log n + 28$ bits and you are more likely to win the Powerball.

$$h(f) = f \pmod{p} \quad \text{for prime } p \in \{2, \dots, tn \log(tn)\}$$

Claim: If $f_1 \neq f_2$ then $h(f_1) = h(f_2)$ with probability just $\frac{2}{t}$.

First observation: If $h(f_1) = h(f_2)$, then:

$$(f_1 - f_2) \pmod{p} = 0.$$

In other words, we only fail if the $f_1 - f_2$ is divisible by p .

Question: What is the chance that $f_1 - f_2$, which is an integer less than 2^n , is divisible by a random prime $p \in \{2, \dots, tn \log(tn)\}$?

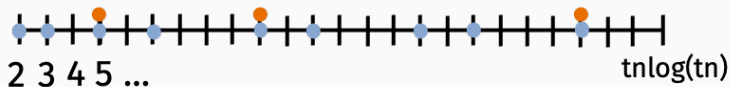
RANDOM FINGERPRINTING

Number of distinct prime factors of $f_1 - f_2$: At most n .

Number of primes between $\{2, \dots, tn \log(tn)\}$: At least $\frac{tn \log(tn)}{\log(tn \log(tn))}$ via prime number theorem.

● = prime number

● = prime factors of $f_1 - f_2$



Chance we pick a prime factor of $f_1 - f_2$ is less than:

$$\frac{n}{\frac{tn \log(tn)}{\log(tn \log(tn))}} = \frac{\log(tn \log(tn))}{t \log(tn)} \leq \frac{2 \log(tn)}{t \log(tn)}$$

Conclusion: The chance that a random prime $p \in \{2, \dots, tn \log(tn)\}$ is a factor of $f_1 - f_2$ is $\leq \frac{2}{t}$.

So, for two files $f_1 \neq f_2$, the chance that $h(f_1) = h(f_2) \leq \frac{2}{t}$.

Set $t = 10^{18}$ (the chance you win the Powerball twice in a row).

Fingerprint size: At most $2 \log_2(nt) = 2 \log(n) + 2 \log_2(10^{18})$ bits.

Suppose we are fingerprinting 1mb image files. $n \approx 8 \cdot 10^6$, so our fingerprint has size:

166 bits

This amounts to a 50,000x reduction over sending and comparing the original files.