

CSE 566 Spring 2023

Universal Hitting Set & Bloom Filters

Instructor: Mingfu Shao


Seeking Orders with Lower Density

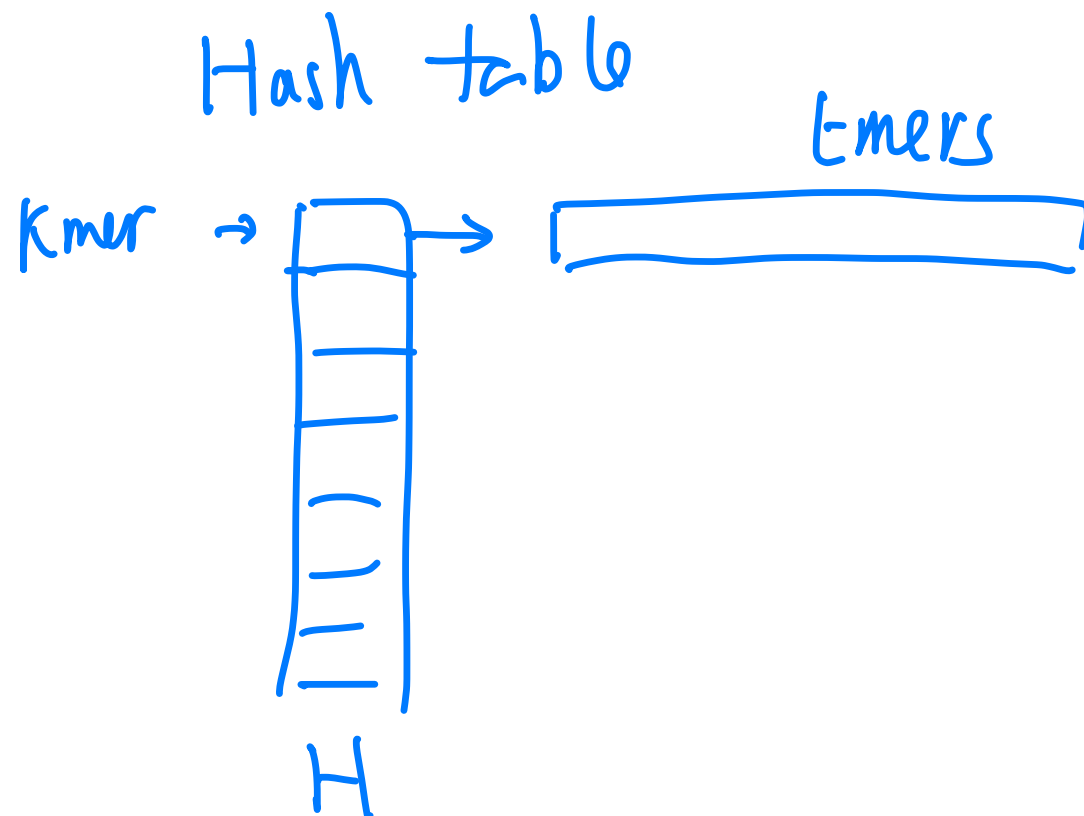
- It is desirable to have a low expected density.
- Open questions:
 - What is the order π with the lowest expected density?
 - Is the lower bound of $1.5/(1 + w)$ reachable?
- Current methods to construct orders with lower expected density are all based on universal hitting set.

Universal Hitting Set (UHS)

- A universal hitting set is a set of kmers H such that every string of length $L = w + k - 1$ has at least one kmer in H .
- A UHS **hits** every string of length L .
- UHS has applications such as indexing sequences. of length L

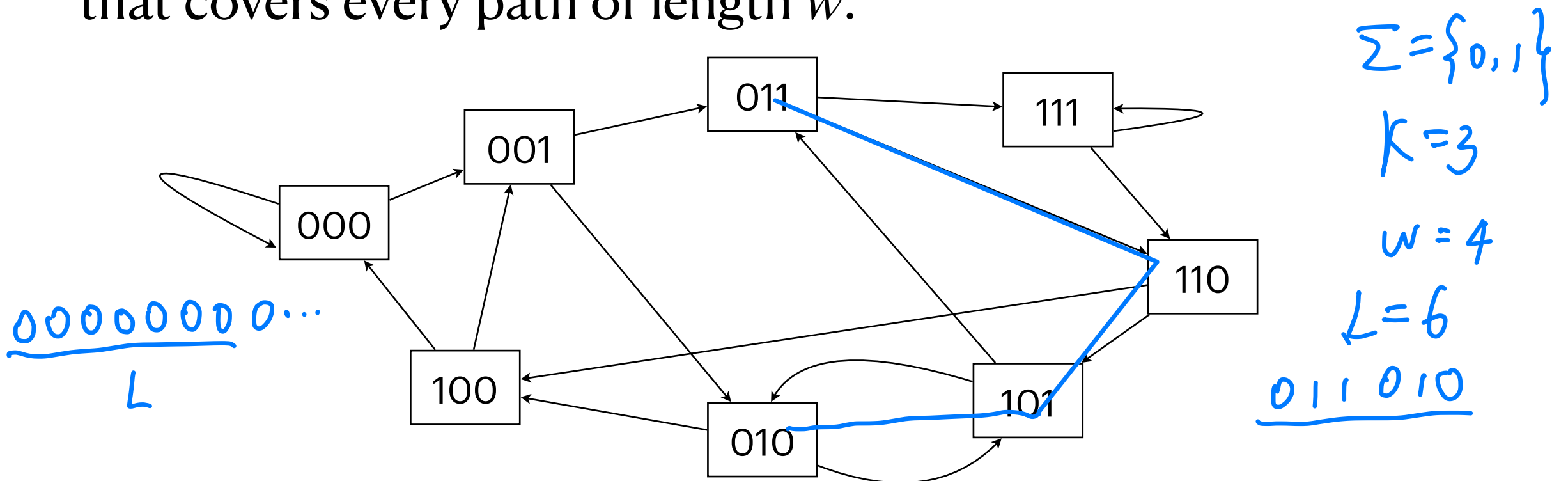
L
 $= k + w - 1$





Finding UHS of Smaller Size

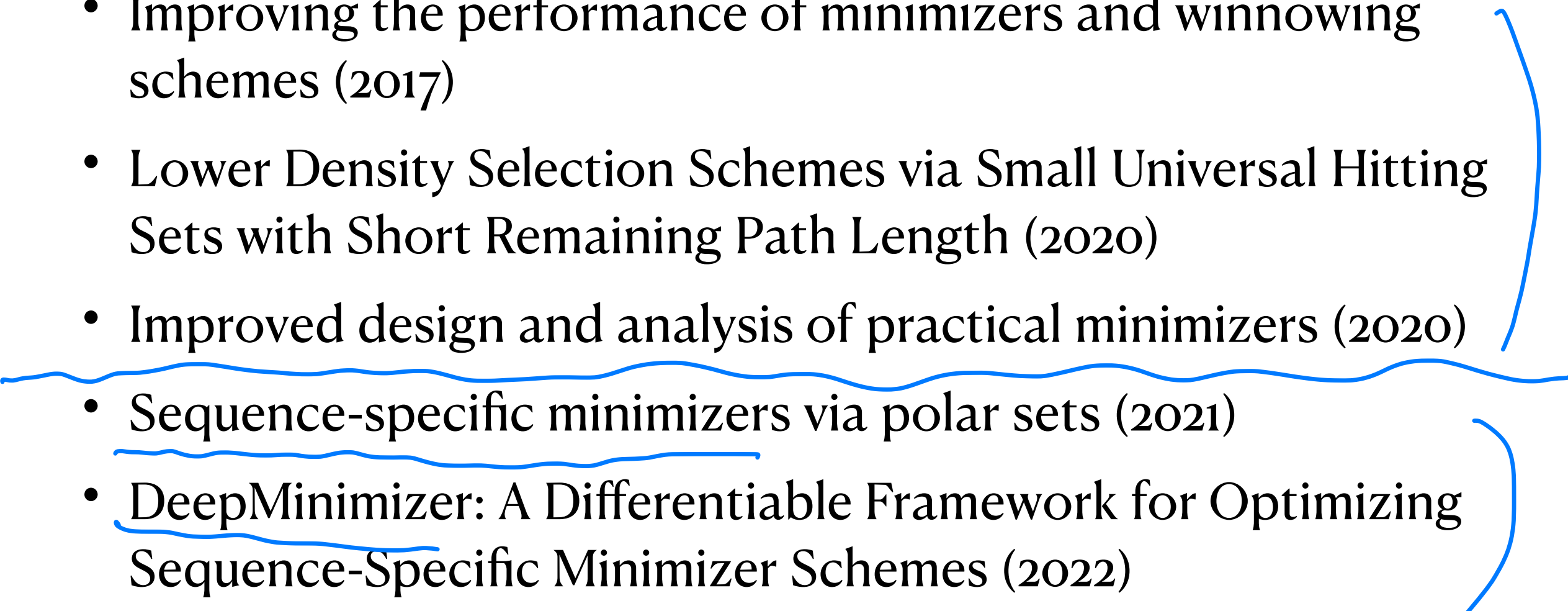
- UHS is closely related to dBG with vertices being Σ^k
- A string of length $L = w + k - 1$ corresponds to a path of length w in the dBG.
- Find a UHS is equivalent to finding a subset of vertices of dBG that covers every path of length w .



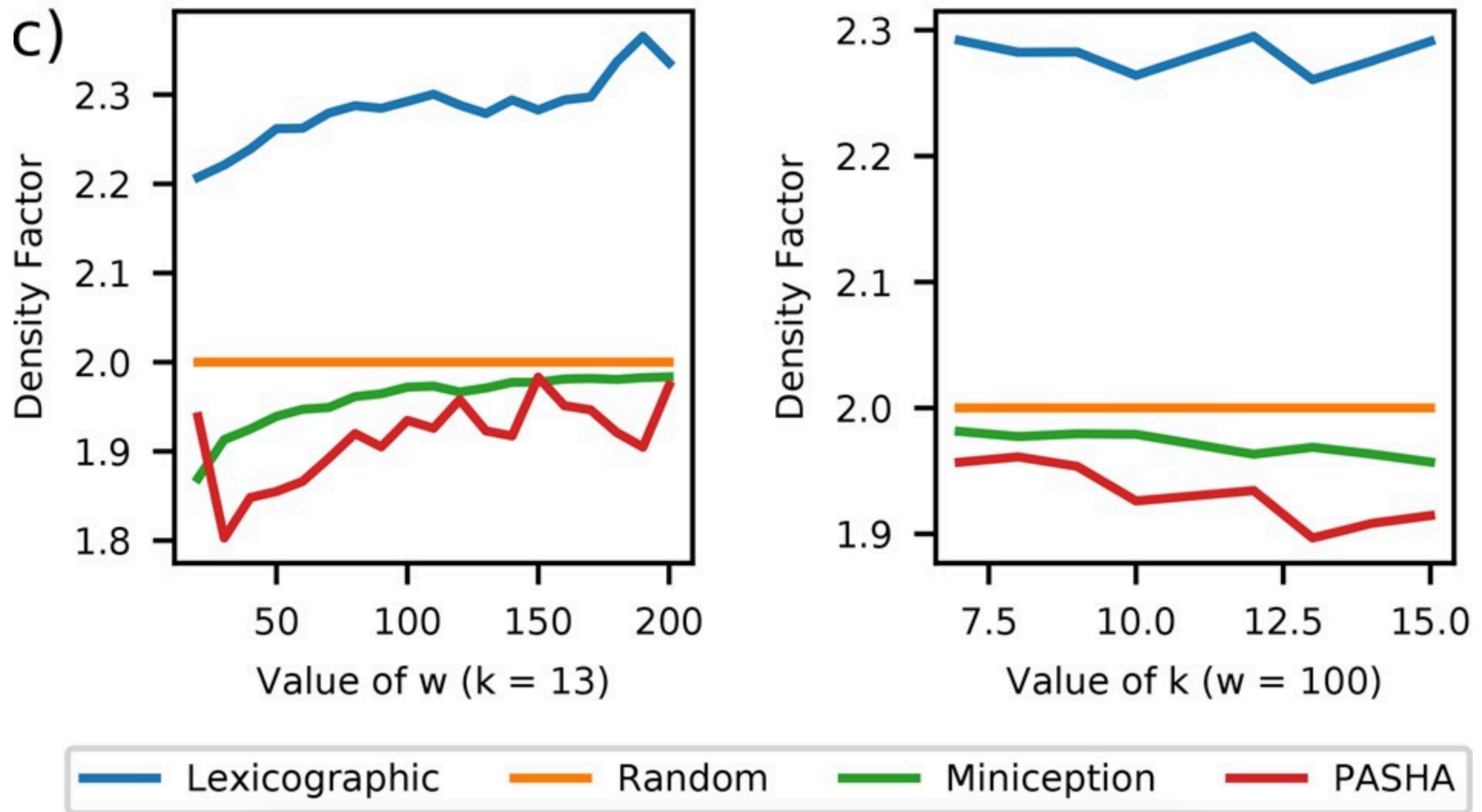
Existing Algorithms to Find UHS

- A greedy approach (DOCKS):
 - To compute vertices that breaks all cycles of dBG;
 - Then to compute vertices that cover all paths of length ℓ . ^{w.}
- Existing work:
 - Compact Universal k-mer Hitting Sets (2016). DOCKS
 - A Randomized Parallel Algorithm for Efficiently Finding Near-Optimal Universal Hitting Sets (2020)
- Remains challenging (scaling only to $k = 16$).

Minimizers with Lower Density

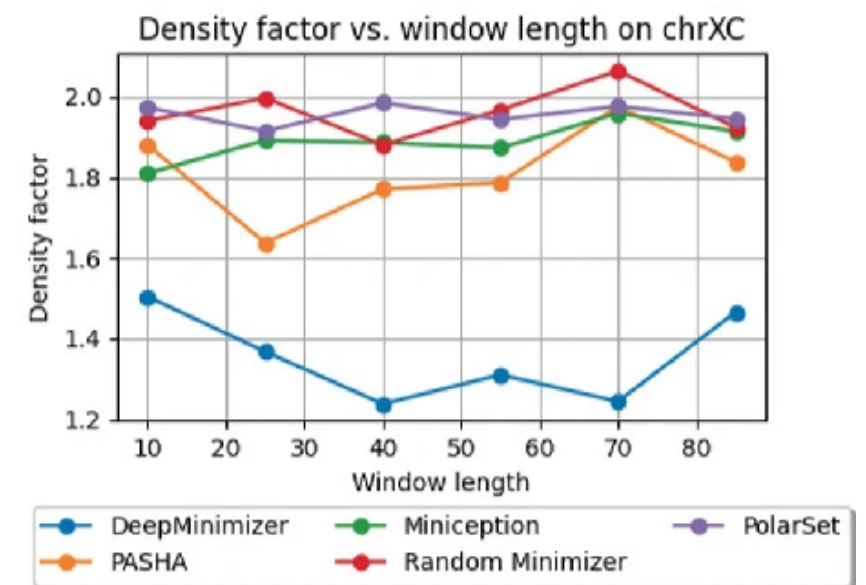
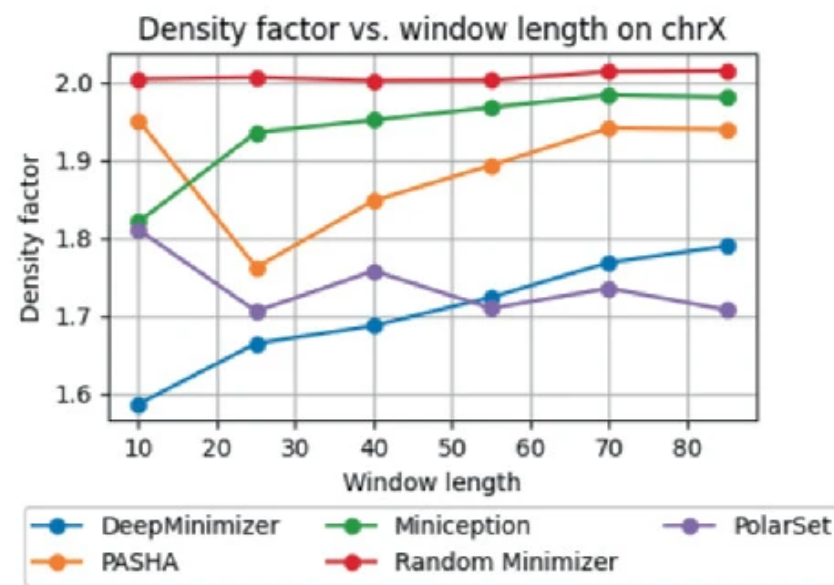
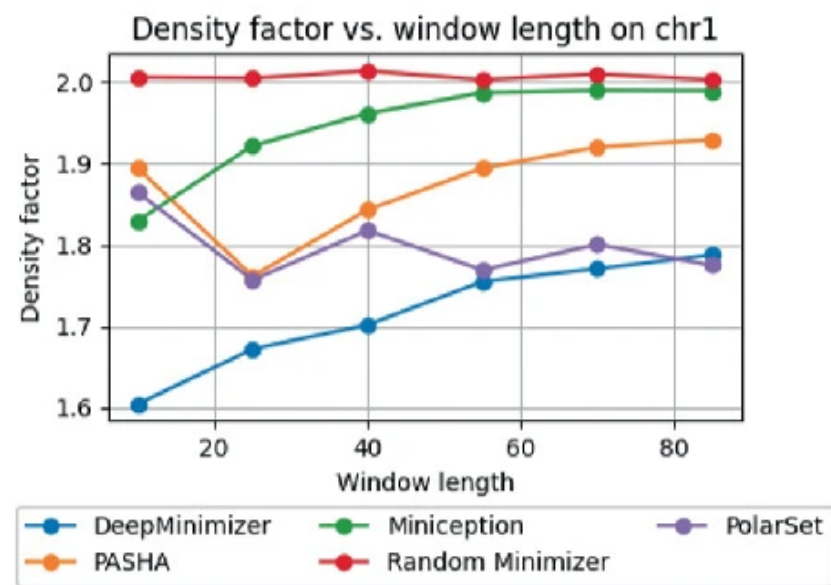
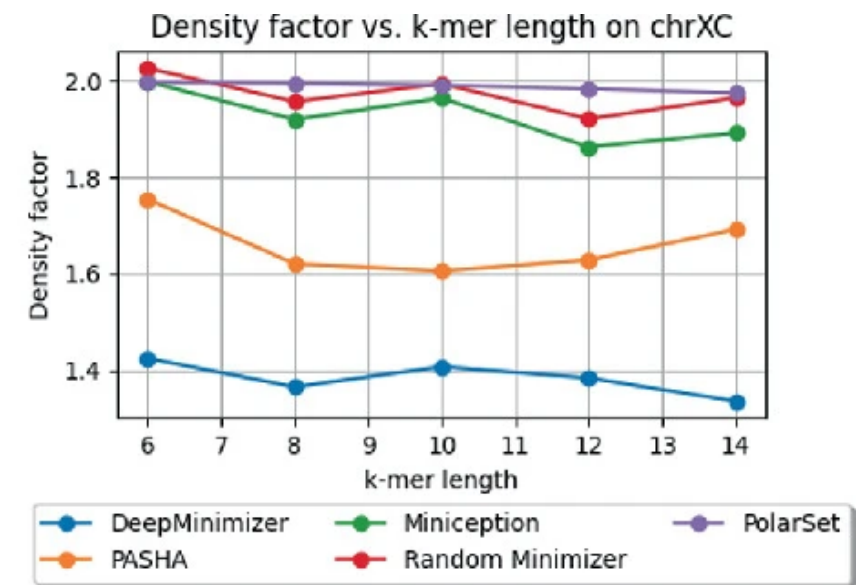
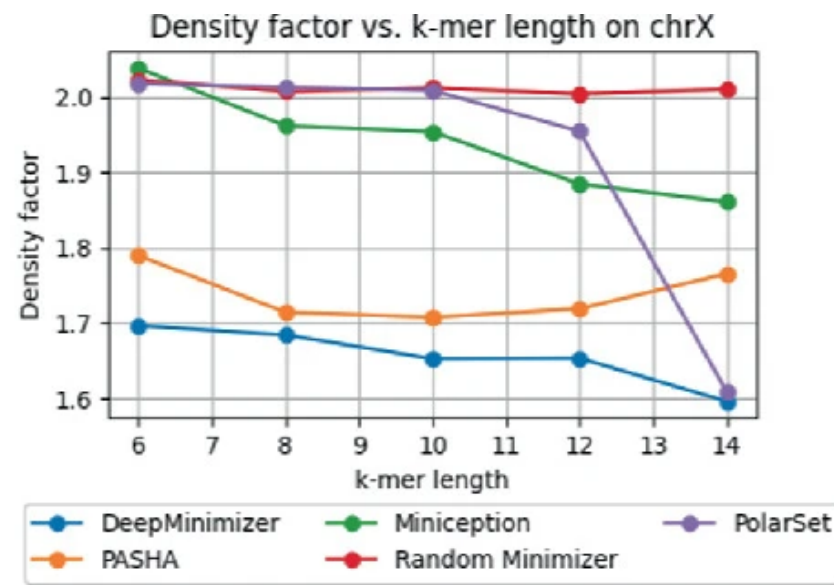
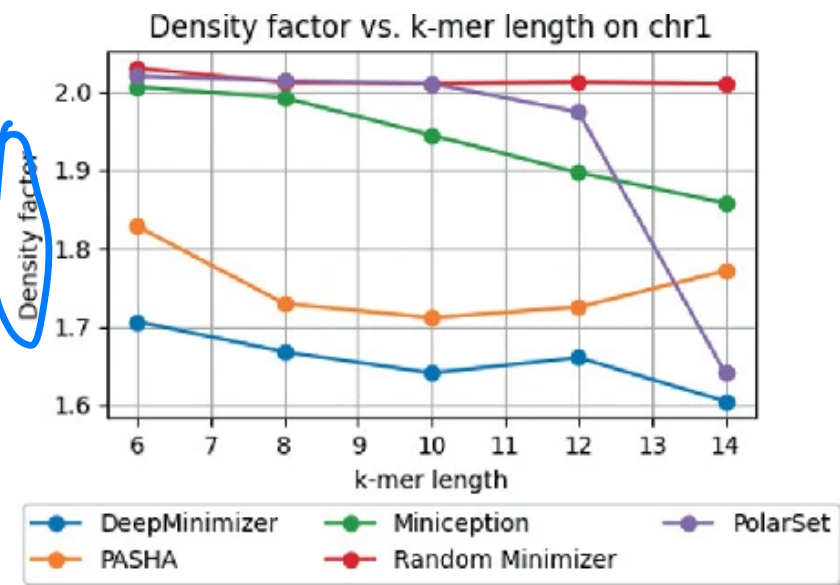
- Improving the performance of minimizers and winnowing schemes (2017)
 - Lower Density Selection Schemes via Small Universal Hitting Sets with Short Remaining Path Length (2020)
 - Improved design and analysis of practical minimizers (2020)
 - Sequence-specific minimizers via polar sets (2021)
 - DeepMinimizer: A Differentiable Framework for Optimizing Sequence-Specific Minimizer Schemes (2022)
- 
- A blue wavy line underlines the last three items of the list. A blue bracket on the right side groups the first three items, and another blue bracket on the right side groups the last two items.

Some Results



Data-Dependent Minimizers

$$\frac{\sigma}{1+w}$$



DeepMinimizer: A Differentiable Framework for Optimizing Sequence-Specific Minimizer Schemes (2022)

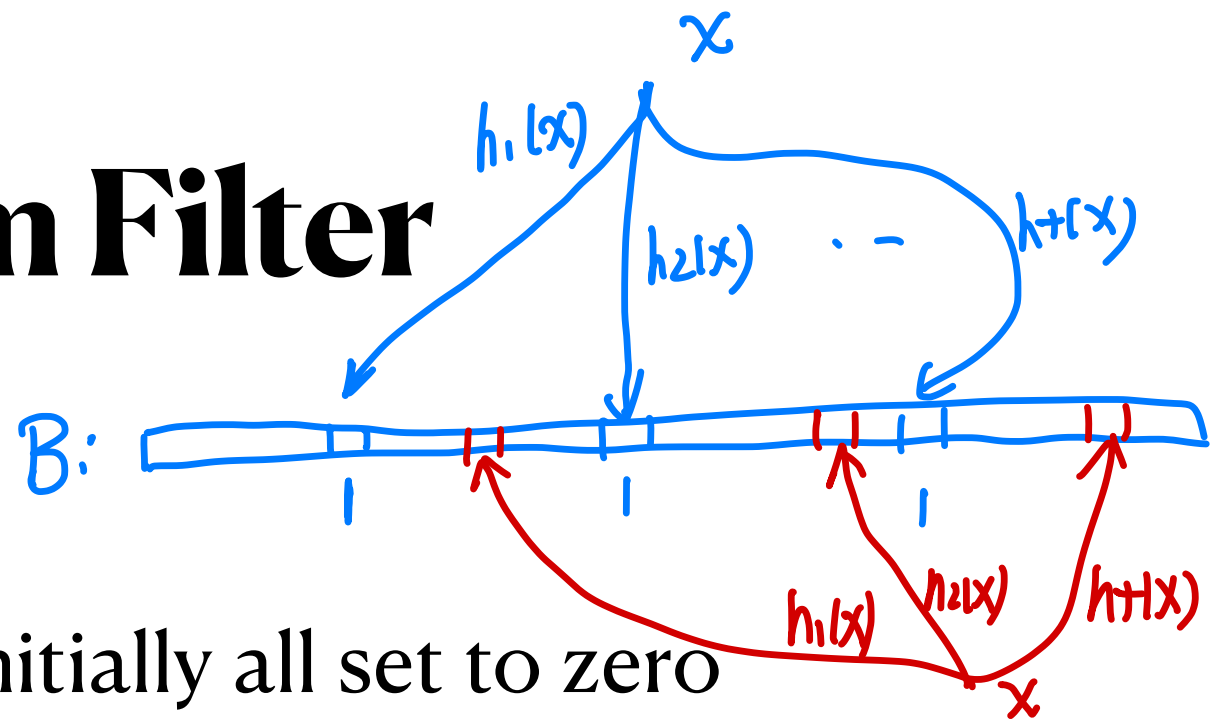
Membership Query

- Store a large set S where memory/space is at a premium
 - All kmers in a large collection of genomic datasets
 - All sentences from a large set of documents
 - All malicious websites/URLs
- A data structure for membership query supports:
 - Add(S, x): add item x to the set S
 - Query(S, x): return true if $x \in S$

Approximate Membership Query


- Two types of errors:
 - False-positive: $\text{Query}(S, x)$ returns true when $x \notin S$
 - False-negative: $\text{Query}(S, x)$ returns false when $x \in S$
- Both are bad but false-negatives are often more severe.
- Bloom filter is a data structure that avoids false-negatives completely but has false-positives.

Bloom Filter



- A bloom filter consists of
 - An bit-array B of size m , initially all set to zero
 - t hash functions h_1, h_2, \dots, h_t , where each h_i maps $x \in U$ to an index of B , i.e., $\{1, 2, \dots, m\}$.
- Add(S, x): set $B[h_i(x)] = 1$, for all $i = 1, 2, \dots, t$
- Query(S, x): return true if $B[h_i(x)] = 1$ for all $i = 1, 2, \dots, t$; otherwise return false.
- A bloom filter allows for querying items without actually storing them!

Minimizing False Positives

- By choosing proper t : t should not be too small or too large
 - **Theorem**: the probability of false positive is minimized when $t = (m/n)\ln 2$, where n is the number of items added.
 - **Proof**: assume that the t hash functions map x to t positions randomly and independently.
- B : 
- After a single item is added, the probability that a particular bit remains zero is: $\Pr = \left(1 - \frac{1}{m}\right)^t$
 - After adding all n items, the probability that a particular bit is set to 1 is: $\Pr = 1 - \left(1 - \frac{1}{m}\right)^{t \cdot n}$

Proof (continued)

- We get a false-positive when all the t bits for a given element x , i.e., $h_1(x), \dots, h_t(x)$, are set to 1; the probability is:

$$\Pr(\text{FP}) = \Pr(h_1(x)=1, h_2(x)=1, \dots, h_t(x)=1)$$

$$= \left(1 - \left(1 - \frac{1}{m}\right)^{tn}\right)^t$$

$$e = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n$$

$$= \left(1 - \boxed{\left(1 - \frac{1}{m}\right)^m}^{\frac{tn}{m}}\right)^t \approx \left(1 - e^{-tn/m}\right)^t.$$

$$\ln \Pr(\text{FP}) = t \cdot \ln \left(1 - e^{-tn/m}\right).$$

Proof (continued)

$$\frac{d \ln \text{Pr}(\text{FP})}{dt} = 1 \cdot \ln(1 - e^{-tn/m}) + t \cdot \frac{e^{-tn/m} \cdot \frac{n}{m}}{1 - e^{-tn/m}}$$

$$\text{let } \frac{d \ln \text{Pr}}{dt} = 0 \Rightarrow \underline{\ln(1 - e^{-tn/m}) - e^{-tn/m} \cdot \ln(1 - e^{-tn/m})}$$
$$= -e^{-tn/m} \cdot \frac{nt}{m} = \underline{e^{-tn/m} \cdot \ln e^{-\frac{nt}{m}}}$$

$$\Rightarrow \ln(1 - e^{-tn/m}) \cdot \underline{(1 - e^{-tn/m})} = \ln e^{-\frac{nt}{m}} \cdot \underline{e^{-tn/m}}$$

$$1 - e^{-tn/m} = e^{-tn/m} \Rightarrow t = \frac{m}{n} \ln 2.$$