

$$1a) \quad 4,000,000 \times 40 \text{ MB} = 160,000,000 \text{ MB}$$

$$1b) \quad B_{k \times m} \rightarrow k \times m = 27,000,000 \text{ bits} \\ = 3.375 \text{ MB}$$

$$1c) \quad FPR = \prod_{i=1}^k P[B_i(h_i(x)) = 1]$$

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

$$k=30 \quad n=4,000,000$$

$$.03 = \left[1 - (1 - \frac{1}{m})^{4,000,000} \right]^{30}$$

$$m = 1,814,524$$

$$1d) = 30 \times 1,814,524 = 54,435,720 \text{ bits}$$

$$6.8 \text{ MB}$$

about double the space but still very short
of if it was deterministic

2a) $.5 \times 102,000 = 61,000$ seconds

2b)

- look for URL in Bloom Filter
- if it returns for malicious
 - check in database to know for sure (this is deterministic)
- if it's not there, do nothing else

2c) $102,000 \times .03 = 3060$ false positives

$$(2000 + 3060)(.5 \text{ seconds}) + (102000 - 3060)(1 \text{ ms})$$

$$2630 \text{ seconds} + 98.94 \text{ seconds}$$

$$2628.94 \text{ seconds}$$

2d) Bloom filter + database deterministic approach is 4% of the run time of the purely database approach

Bloom filter is much better approach

Question 2)

Because there is only one row the

$$FPR = [1 - (1 - \frac{1}{m})^n]$$

instead of

$$FPR = [1 - (1 - \frac{1}{m})^n]^k \quad \text{because there}$$

are no longer k rows in the matrix

b) if $FPR = .05$

$$k = 30 \quad n = 4,000,000$$

$$.05 = [1 - (1 - 1/m)^n]$$

$$.95 = (1 - 1/m)^{4,000,000}$$

$$\log(.95) = 4,000,000 \log(1 - 1/m)$$

$$-1.28 \times 10^{-8} = \log(1 - 1/m)$$

$$-1.28 \times 10^{-8} = \log(e^{-1/m})$$

$$-1.28 \times 10^{-8} = -1/m$$

$$m = 77,982,909 \text{ bits}$$

$$\text{space} = 9.747 \text{ MB}$$

modified version of
bloom filter is
worse

$$.05 = [1 - (1 - 1/m)^n]^k$$

$$.05 = [1 - (1 - 1/m)^{4,000,000}]^{30}$$

$$\frac{\log(.05)}{30} = \log(1 - (1 - 1/m)^{4,000,000})$$

$$-.099 = \log(1 - (1 - 1/m)^{4,000,000})$$

$$e^{-.099} = (1 - (1 - 1/m)^{4,000,000})$$

$$.09 = (1 - 1/m)^{4,000,000}$$

$$\frac{\log(.09)}{4,000,000} = \log(e^{-1/m})$$

$$-6.0198 \times 10^{-7} = -1/m$$

$$m = 1,661,184$$

$$k = 30$$

$$\text{total space} = m \times k$$

$$= 6.22 \text{ MB}$$

Coding Challenge

$$\text{original FPR} = .0359$$

$$\text{FPR with } k=10 \text{ } m=8000 \text{ and } t=(1, m-1) = 1.0$$

$$\text{FPR with } k=10 \text{ } m=79,000 \text{ and } t(1, m-1) = .0339$$

$$\text{ratio from 2b} = \frac{77,982,909}{1,661,182} = 46.944$$

$$\text{ratio from coding challenge} = \frac{79,000}{8000} = 9.875$$

I think the ratio for the homework is smaller because $n=10,000$ but for the derivation $n=4,000,000$, k was also larger for the derivation

