# APPENDIX U
# MATHEMATICAL BASIS OF THE BIRTHDAY ATTACK

## William Stallings

Copyright 2013

In this appendix, we derive the mathematical justification for the birthday attack. We begin with a related problem and then look at the problem from which the name "birthday attack" is derived.

## U.1  RELATED PROBLEM

A general problem relating to hash functions is the following. Given a hash function H, with $n$ possible outputs and a specific value H($x$), if H is applied to $k$ random inputs, what must be the value of $k$ so that the probability that at least one input $y$ satisfies H($y$) = H($x$) is 0.5?

For a single value of $y$, the probability that H($y$) = H($x$) is just $1/n$. Conversely, the probability that H($y$) ≠ H($x$) is [1 – (1/n)]. If we generate $k$ random values of $y$, then the probability that none of them match is just the product of the probabilities that each individual value does not match, or [1 – (1/n)]$^k$. Thus, the probability that there is at least one match is 1 – [1 – (1/n)]$^k$.

The binomial theorem can be stated as

$$\left(1-a\right)^{k} = 1 - ka + \frac{k\left(k-1\right)}{2!}a^{2} - \frac{k\left(k-1\right)\left(k-2\right)}{3!}a^{3} \ldots$$

For very small values of $a$, this can be approximated as (1 – $ka$). Thus, the probability of at least one match is approximated as 1 – [1 – (1/n)]$^k$ ≈ 1 – [1 – (k/n)] = k/n. For a probability of 0.5, we have $k = n/2$.

In particular, for an $m$-bit hash code, the number of possible codes is $2^m$ and the value of $k$ that produces a probability of one-half is

$$k = 2^{(m-1)} \tag{1}$$

## U.2  THE BIRTHDAY PARADOX

The birthday paradox is often presented in elementary probability courses to demonstrate that probability results are sometimes counterintuitive. The problem can be stated as follows: What is the minimum value of $k$ such that the probability is greater than 0.5 that at least two people in a group of $k$ people have the same birthday? Ignore February 29 and assume that each birthday is equally likely.

We can reason to the answer as follows. The probability that the birthdays of any two people are not alike is clearly 364/365 (since there is only one chance in 365 that one person's birthday will coincide with another's). The probability that a third person's birthday will differ from the other two is 363/365; a fourth person's, 362/365; and so on, until we reach the 24th person (342/365). We thus obtain a series of 23 fractions which must be multiplied together to reach the probability that all 24 birthdays are different. The product is a fraction that reduces to about 0.507, or slightly better than 1/2, for a coincidence among 23 people.

To derive this answer formally, let us define

$P(n, k)$  =  Pr[at least one duplicate in $k$ items, with each item able to take on one of $n$ equally likely values between 1 and $n$]

Thus, we are looking for the smallest value of $k$ such that $P(365, k) \geq 0.5$. It is easier first to derive the probability that there are no duplicates, which we designate as $Q(365, k)$. If $k > 365$, then it is impossible for all values to be different. So we assume $k \leq 365$. Now consider the number of different ways, $N$, that we can have $k$ values with no duplicates. We may choose any of the 365 values for the first item, any of the remaining 364

numbers for the second item, and so on. Hence, the number of different ways is

$$N = 365 \times 364 \times \ldots \left(365 - k + 1\right) = \frac{365!}{\left(365 - k\right)!}$$  **(2)**

If we remove the restriction that there are no duplicates, then each item can be any of 365 values, and the total number of possibilities is $365^k$. So the probability of no duplicates is simply the fraction of sets of values that have no duplicates out of all possible sets of values:

$$Q\left(365, k\right) = \frac{365! / \left(365 - k\right)!}{\left(365\right)^k} = \frac{365!}{\left(365 - k\right)! \left(365\right)^k}$$

and

$$P\left(365, k\right) = 1 - Q\left(365, k\right) = 1 - \frac{365!}{\left(365 - k\right)! \left(365\right)^k}$$  **(3)**

This function is plotted in Figure U.1. The probabilities may seem surprisingly large to anyone who has not considered the problem before. Many people would guess that to have a probability greater than 0.5 that there is at least one duplicate, the number of people in the group would have to be about 100. In fact, the number is 23, with P(365, 23) = 0.5073. For $k$ = 100, the probability of at least one duplicate is 0.9999997.
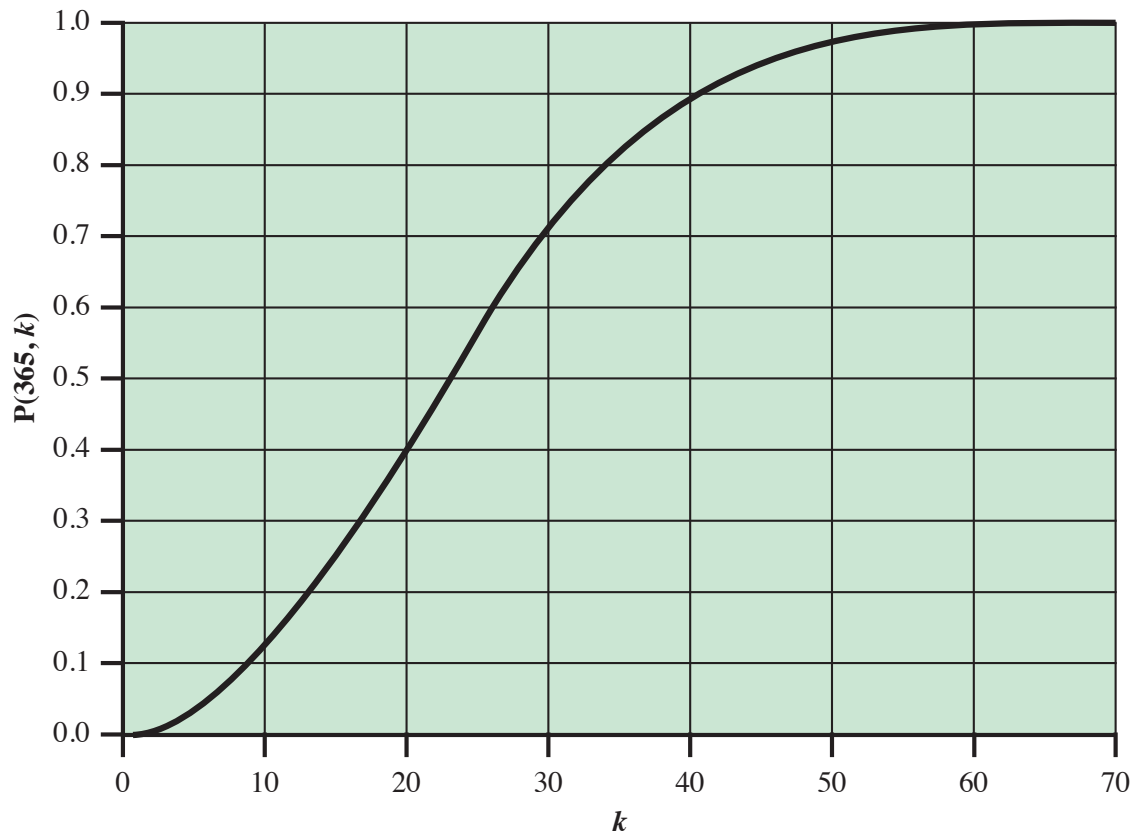
**Figure U.1 The Birthday Paradox**

Perhaps the reason that the result seems so surprising is that if you consider a particular person in a group, the probability that some other person in the group has the same birthday is small. But the probability that we are concerned with is the probability that *any* pair of people in the group has the same birthday. In a group of 23, there are $(23(23 - 1))/2 = 253$ different pairs of people. Hence the high probabilities.

## U.3  USEFUL INEQUALITY

Before developing a generalization of the birthday problem, we derive an inequality that will be needed:

$$(1 - x) \leq e^{-x} \qquad \text{for all } x \geq 0 \qquad \textbf{(4)}$$
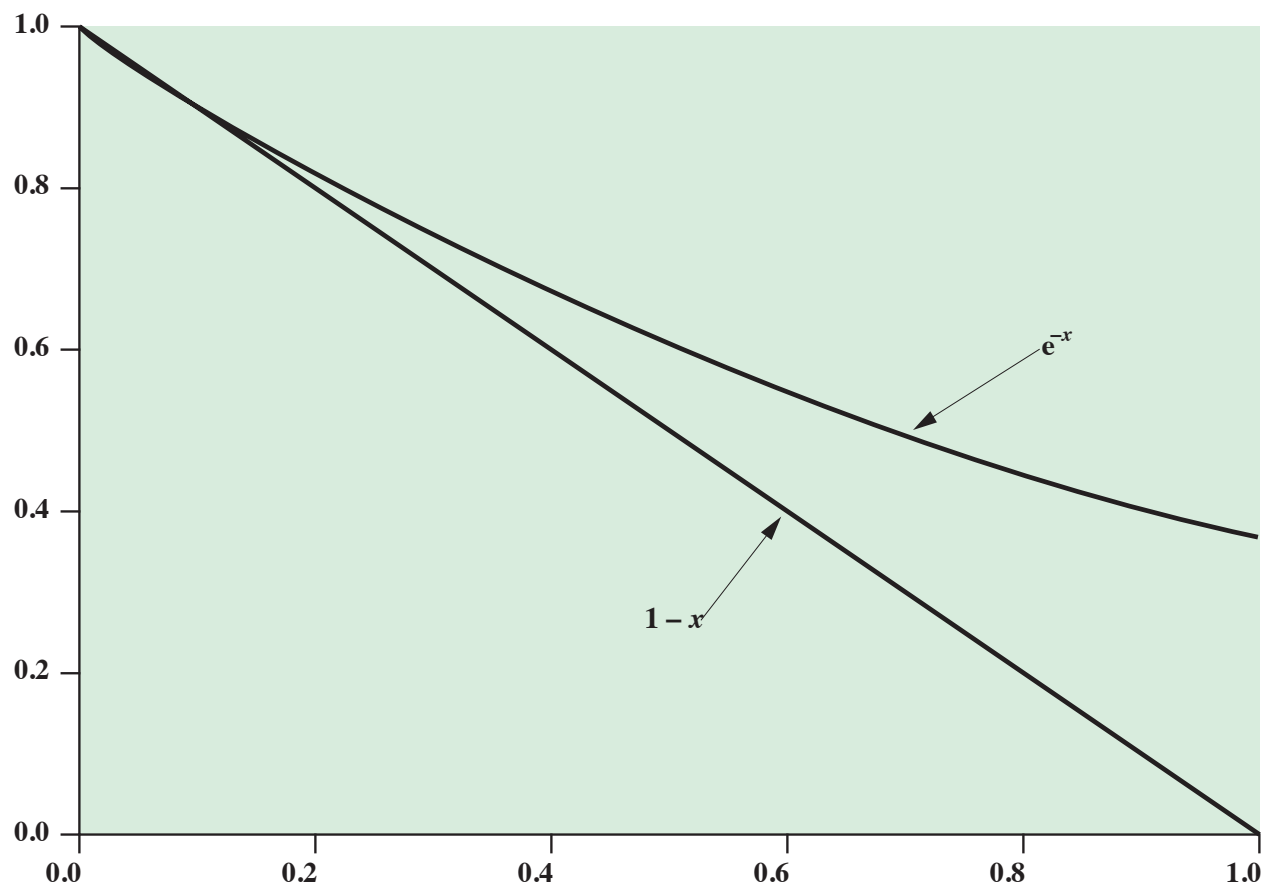
**Figure U.2  A Useful Inequality**

Figure U.2 illustrates the inequality. To see that the inequality holds, note that the lower line is the tangent to $e^{-x}$ at $x = 0$. The slope of that line is just the derivative of $e^{-x}$ at $x = 0$:

$$f(x) = e^{-x}$$
$$f'(x) = \frac{d}{dx}e^{-x} = -e^{-x}$$
$$f'(0) = -1$$

U-6

The tangent is a straight line of the form $ax + b$ with $a = -1$, and the tangent at $x = 0$ must equal $e^{-0} = 1$. Thus, the tangent is the function $(1 - x)$, confirming the inequality of Equation (11.4). Further, note that for small $x$, we have $(1 - x) \approx e^{-x}$.

## U.4  THE GENERAL CASE OF DUPLICATIONS

The birthday problem can be generalized to the following problem. Given a random variable that is an integer with uniform distribution between 1 and $n$ and a selection of $k$ instances ( $k \leq n$) of the random variable, what is the probability, $P(n, k)$, that there is at least one duplicate? The birthday problem is just the special case with $n = 365$. By the same reasoning as before, we have the following generalization of Equation (3):

$$P\left(n,k\right) = 1 - \frac{n!}{\left(n-k\right)! n^k} \tag{5}$$

We can rewrite this as

$$
\begin{aligned}
P\left(n,k\right) &= 1 - \frac{n \times \left(n-1\right) \times \ldots \times \left(n-k+1\right)}{n^k} \\
&= 1 - \left[\frac{n-1}{n} \times \frac{n-2}{n} \times \ldots \times \frac{n-k+1}{n}\right] \\
&= 1 - \left[\left(1 - \frac{1}{n}\right) \times \left(1 - \frac{2}{n}\right) \times \ldots \times \left(1 - \frac{k-1}{n}\right)\right]
\end{aligned}
$$

Using the inequality of Equation (4),

$$P(n,k) > 1 - \left[ \left(e^{-1/n}\right) \times \left(e^{-2/n}\right) \times \ldots \times \left(e^{-(k-1)/n}\right) \right]$$

$$> 1 - e^{-\left[(1/n)+(2/n)+\ldots+((k-1)/n)\right]}$$

$$> 1 - e^{-(k\times(k-1))/2n}$$

Now let us pose the question: What value of $k$ is required such that $P(n, k)$ > 0.5? To satisfy the requirement, we have

$$1/2 = 1 - e^{-(k\times(k-1))/2n}$$

$$2 = e^{(k\times(k-1))/2n}$$

$$\ln 2 = \frac{k \times (k-1)}{2n}$$

For large $k$, we can replace $k \times (k-1)$ by $k^2$, and we get

$$k = \sqrt{2(\ln 2)n} = 1.18\sqrt{n} \approx \sqrt{n} \tag{6}$$

As a reality check, for $n = 365$, we get $k = 1.18 \times \sqrt{365} = 22.54$, which is very close to the correct answer of 23.

We can now state the basis of the birthday attack in the following terms. Suppose we have a function H, with $2^m$ possible outputs (i.e., an $m$-bit output). If H is applied to $k$ random inputs, what must be the value of $k$ so that there is the probability of at least one duplicate [i.e., H($x$) = H($y$) for some inputs $x$, $y$)]? Using the approximation in Equation (6),

$$k = \sqrt{2^m} = 2^{m/2} \tag{7}$$

## U.5  OVERLAP BETWEEN TWO SETS

There is a problem related to the general case of duplications that is also of relevance for our discussions. The problem is this: Given an integer random variable with uniform distribution between 1 and $n$ and two sets of $k$ instances ( $k \leq n$) of the random variable, what is the probability, R($n$, $k$), that the two sets are not disjoint; that is, what is the probability that there is at least one value found in both sets?

Let us call the two sets X and Y, with elements $\{x_1, x_2, \ldots, x_k\}$ and $\{y_1, y_2, \ldots, y_k\}$, respectively. Given the value of $x_1$, the probability that $y_1 = x_1$ is just $1/n$, and therefore the probability that $y_1$ does not match $x_1$ is [1 − (1/n)]. If we generate the $k$ random values in Y, the probability that none of these values is equal to $x_1$ is $[1 − (1/n)]^k$. Thus, the probability that there is at least one match to $x_1$ is $1 − [1 − (1/n)]^k$.

To proceed, let us make the assumption that all the elements of X are distinct. If $n$ is large and if $k$ is also large (e.g., on the order of $\sqrt{n}$ ), then this is a good approximation. In fact, there may be a few duplications, but most of the values will be distinct. With that assumption, we can make the following derivation:

$$\Pr[\text{no match in } Y \text{ to } x_1] = \left(1 - \frac{1}{n}\right)^k$$

$$\Pr[\text{no match in } Y \text{ to } X] = \left(\left(1 - \frac{1}{n}\right)^k\right)^k = \left(1 - \frac{1}{n}\right)^{k^2}$$

$$\mathrm{R}(n,k) = \Pr[\text{at least one match in } Y \text{ to } X] = 1 - \left(1 - \frac{1}{n}\right)^{k^2}$$

Using the inequality of Equation (4),

$$R(n,k) > 1 - \left(e^{-1/n}\right)^{k^2}$$

$$R(n,k) > 1 - \left(e^{-k^2/n}\right)$$

Let us pose the question: What value of $k$ is required such that R($n$, $k$) > 0.5? To satisfy the requirement, we have

$$1/2 = 1 - \left(e^{-k^2/n}\right)$$

$$2 = e^{k^2/n}$$

$$\ln(2) = \frac{k^2}{n} \qquad \textbf{(8)}$$

$$k = \sqrt{\left(\ln(2)\right)n} = 0.83\sqrt{n} \approx \sqrt{n}$$

We can state this in terms related to birthday attacks as follows. Suppose we have a function H, with $2^m$ possible outputs (i.e., an $m$-bit output). Apply H to $k$ random inputs to produce the set X and again to $k$ additional random inputs to produce the set Y. What must be the value of $k$ so that there is the probability of at least 0.5 that there is a match between the two sets (i.e., H($x$) = H($y$) for some inputs $x \in$ X, $y \in$ Y)? Using the approximation in Equation (8):

$$k = \sqrt{2^m} = 2^{m/2}$$