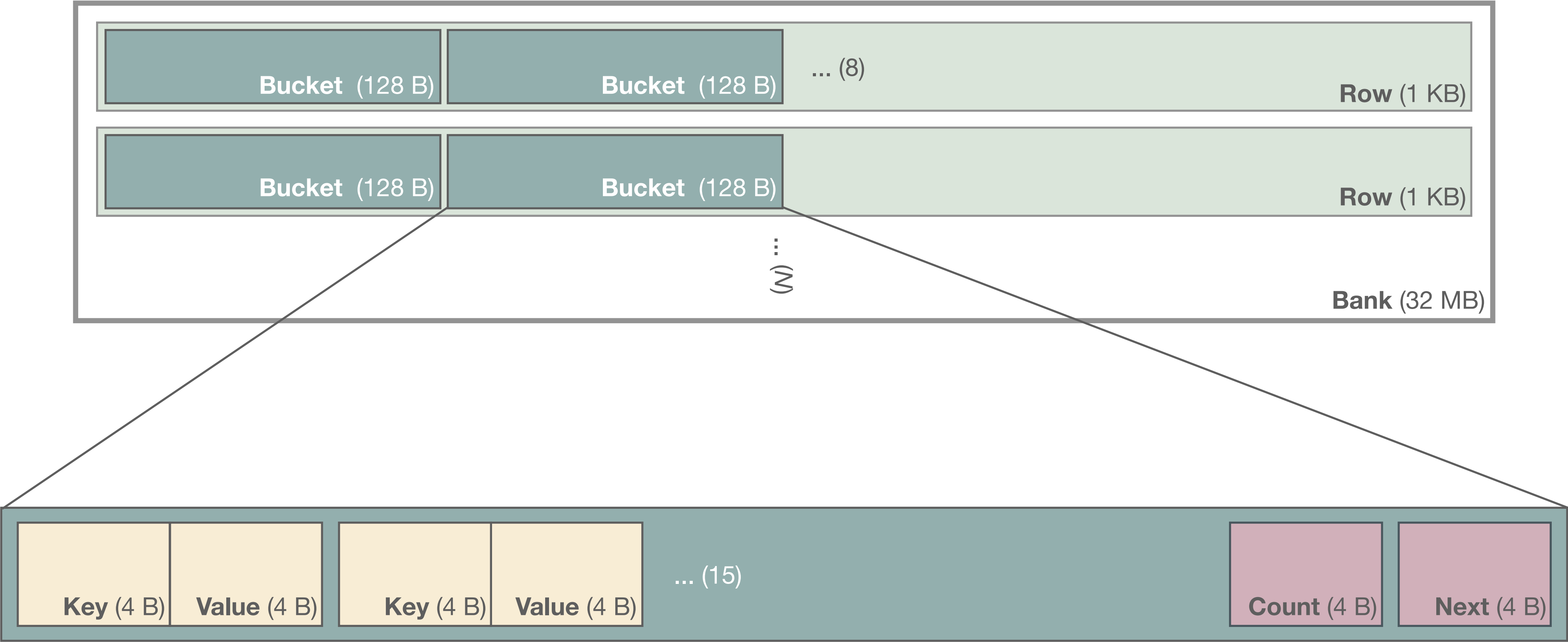# Intra-bank hash maps in BLIMP

**Kevin Gaffney | February 17, 2023**

# Motivation

- We need an efficient data structure that enables random lookups to perform hash joins and hash aggregates.

- We need to tailor the design to BLIMP's unique characteristics.

# Overview

| **Insert** |
| --- |
| **Input**: a (key, value) pair |
| ```
bucket = buckets[hash(key)];

WHILE (bucket.next != NULL) {

    bucket = bucket.next;

}

IF (bucket.count == 15) {

    bucket = new_bucket(bucket);

}

bucket[bucket.count] = (key, value);

++bucket.count;
``` |

**Get**

**Input**: a key.

**Output**: a pointer to a value (may be NULL).

```
bucket = buckets[hash(key)];

WHILE (bucket != NULL) {

    FOR (i IN [0, bucket.count)) {

        IF (bucket[i].key == key) {

            RETURN pointer to bucket[i].value;

        }

    }

    bucket = bucket.next;

}

RETURN NULL;
```

# Modeling

For 1 million elements and varying load factor $\alpha$, what is the probability $P$ that we find a given item in the **first** bucket we check?

The probability $P$ is given by

$$P(\alpha) = \sum_{i=0}^{\infty} F\left( i,\ 10^6,\ \alpha\frac{15}{10^6} \right) \cdot \frac{15}{\max(i,15)}$$

Probability of a super-bucket having $i$ elements

Probability of finding an element in the first bucket of a super-bucket.

where $F(k, n, p)$ is the probability mass function for the binomial distribution with number of trials $n$ and probability of success $p$.

# Modeling

For 1 million elements and varying load factor $\alpha$, what is the probability $P$ that we find a given item in the **first** bucket we check?