

HW #5

6.41

Given:

- 1) $\text{sizeof(char)} = 1$
- 2) $\text{sizeof(int)} = 4$
- 3) buffer begins at memory address 0
- 4) cache is initially empty
- 5) memory accesses are to the entries of the array buffer.
- 6) cache has 4-byte lines
- 7) 640×480 array of pixels (2-d array)

$$\text{Miss Rate} = 1 - \text{Hit Rate}$$

Calculating Hit Rate

First, we must calculate the size of our pixel struct. Since we have four char elements in our struct, we know that our struct is 4 bytes long.

Next, we visualize the pixel buffer below. We know that our pixel buffer is a 2-dimensional array of size 480×640 .

Thus, the array appears as follows, where each p_i is 4 bytes:

$p_{0,0}, p_{0,1}, p_{0,2}, \dots, p_{0,479}, p_{1,0}, p_{1,1}, p_{1,2}, \dots, p_{1,479}, \dots, p_{479,0}, p_{479,1}, p_{479,2}, \dots, p_{479,479}$

$p_{0,0}, p_{0,1}, p_{0,2}, \dots, p_{0,479}, p_{1,0}, p_{1,1}, p_{1,2}, \dots, p_{1,479}, \dots, p_{479,0}, p_{479,1}, p_{479,2}, \dots, p_{479,479}$

where each $p_i = \begin{bmatrix} r & g & b & a \end{bmatrix}$

Since we know that the cache is initially empty, we know that there will be a miss whenever we will read the red in the struct. However, the subsequent 3 accesses of green, blue, and a will be hits because of spatial locality and the fact that arrays are contiguous stored in memory, as well as the cache is 4-bytes long. This pattern continues at a constant rate for every p_i . Therefore, the hit rate of our cache simplifies to $3/4$.

Calculating Hit Rate

$$\text{Miss Rate} = 1 - \text{Hit Rate}$$

$$= 1 - \frac{3}{4} = \frac{1}{4}$$

6.45

One optimization we could perform on the transpose function is loop unrolling. The idea behind loop unrolling is the fact that we should reduce loop iterations by increasing work per iteration. The effect of loop unrolling is that it makes our program run faster because it reduces the number of iterations and enables the program to exploit instruction-level parallelism. In this case, the work being done would be the striding of columns and rows. Therefore, our program becomes:

```
void transpose(int *dst, int *src, int dim) {
    int i, j;
    for (i = 0; i < dim; i++) {
        for (j = 0; j < dim/4; j++) {
            dst[(j)*dim + i] = src[i*dim + j];
            dst[(j+1)*dim + i] = src[i*dim + (j+1)];
            dst[(j+2)*dim + i] = src[i*dim + (j+2)];
            dst[(j+3)*dim + i] = src[i*dim + (j+3)];
        }
    }
}
```

Another optimization we could make is tiling. Tiling would exploit the fact that the transpose is inherently access within a temporal local setting. Therefore, we could split our tiles into 8×8 tiles in order to exploit the benefit of tiling.

```
void transpose(int *dst, int *src, int dim) {
    int i, j;
    int tile = (dim / 8);
    for (i = 0; i < dim; i++) {
        for (j = 0; j < tile; j++) {
            dst[(j)*dim + i] = src[i*dim + (j)];
            dst[(j+1)*dim + i] = src[i*dim + (j+1)];
            dst[(j+2)*dim + i] = src[i*dim + (j+2)];
        }
    }
}
```


$$\begin{aligned}
 \text{dst}[(j+3)^* \text{dim} + i] &= \text{src}[i^* \text{dim} + (j+3)]; \\
 \text{dst}[(j+4)^* \text{dim} + i] &= \text{src}[i^* \text{dim} + (j+4)]; \\
 \text{dst}[(j+5)^* \text{dim} + i] &= \text{src}[i^* \text{dim} + (j+5)]; \\
 \text{dst}[(j+6)^* \text{dim} + i] &= \text{src}[i^* \text{dim} + (j+6)]; \\
 \text{dst}[(j+7)^* \text{dim} + i] &= \text{src}[i^* \text{dim} + (j+7)];
 \end{aligned}$$

}

}

}