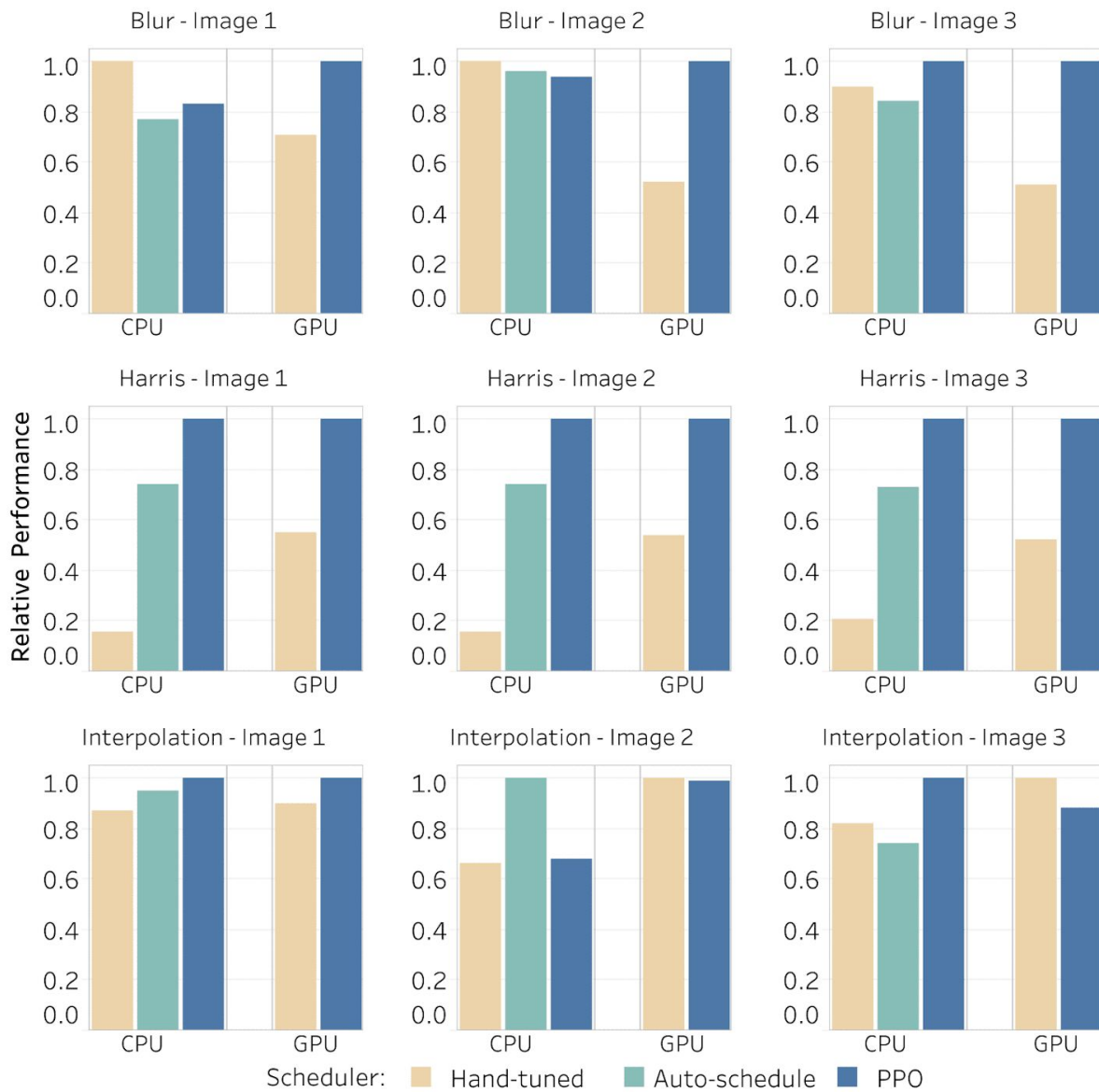# Optimization of Halide Image Processing Schedules with Reinforcement Learning

## --- Generated Schedules Comparison ---

1. **Performance relative to the best result by architecture and scheduling method. The bigger the better. The relative performance is based on the ratio between the execution times of the compared schedules. The best has value 1.0 and the others have lower values.**

**2. Absolute Execution Time (ms) and Relative Slowdown (x) by Architecture and Scheduling Method.**

| | | CPU | | | | | | GPU | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Hand-tuned | | Auto-sched. | | PPO | | Hand-tuned | | PPO | |
| | | $ms$ | $\times$ | $ms$ | $\times$ | $ms$ | $\times$ | $ms$ | $\times$ | $ms$ | $\times$ |
| Blur | Img1 | **0.10** | - | 0.13 | 1.3 | 0.12 | 1.2 | 0.07 | 1.4 | **0.05** | - |
| | Img2 | **0.49** | - | 0.51 | 1.0 | 0.52 | 1.1 | 0.21 | 1.9 | **0.11** | - |
| | Img3 | 2.94 | 1.1 | 3.15 | 1.2 | **2.66** | - | 0.76 | 1.9 | **0.39** | - |
| Harris | Img1 | 3.21 | 6.4 | 0.68 | 1.4 | **0.50** | - | 0.22 | 1.8 | **0.12** | - |
| | Img2 | 13.07 | 6.2 | 2.82 | 1.3 | **2.10** | - | 0.72 | 1.8 | **0.39** | - |
| | Img3 | 36.89 | 4.7 | 10.71 | 1.4 | **7.78** | - | 2.79 | 1.9 | **1.46** | - |
| Interp. | Img1 | 4.51 | 1.2 | 4.10 | 1.0 | **3.91** | - | 3.27 | 1.1 | **2.94** | - |
| | Img2 | 17.43 | 1.5 | **11.49** | - | 16.88 | 1.5 | **6.22** | - | 6.26 | 1.0 |
| | Img3 | 77.23 | 1.2 | 85.89 | 1.4 | **63.57** | - | **16.23** | - | 18.52 | 1.1 |

**3. Best PPO Generated Schedule (CPU) for each Reinforcement Learning Trial:**

| | | |
|---|---|---|
| Blur | Trial 1 | blur_y.tile(x, y, xi, yi, 512, 16);<br>blur_y.parallel(y);<br>blur_x.compute_at(blur_y, yi);<br>blur_x.compute_at(blur_y, x); |
| | Trial 2 | blur_y.unroll(x, 4);<br>blur_y.parallel(y);<br>blur_y.bound(x, 0, input.width()); |
| | Trial 3 | blur_y.tile(x, y, xi, yi, 128, 8);<br>blur_y.unroll(x, 2);<br>blur_x.compute_at(blur_y, x);<br>blur_y.parallel(y);<br>blur_y.bound(y, 0, input.height());<br>blur_y.unroll(x, 4); |
| | Trial 4 | blur_y.tile(x, y, xi, yi, 512, 16);<br>blur_y.parallel(y);<br>blur_x.compute_at(blur_y, x); |
| Harris | Trial 1 | shifted.tile(x, y, xi, yi, 256, 64);<br>shifted.parallel(y);<br>gray.compute_at(shifted, x);<br>ly.compute_at(shifted, x);<br>lx.compute_at(shifted, x); |

| | | |
|---|---|---|
| | Trial 2 | shifted.tile(x, y, xi, yi, 512, 8);<br>shifted.parallel(y);<br>gray.compute_at(shifted, x);<br>ly.compute_at(shifted, x);<br>lx.compute_at(shifted, x);<br>gray.store_at(shifted, y);<br>gray.unroll(x, 2);<br>gray.unroll(x, 4); |
| | Trial 3 | shifted.tile(x, y, xi, yi, 512, 32);<br>lx.compute_at(shifted, x);<br>shifted.parallel(y);<br>gray.compute_at(shifted, x);<br>ly.compute_at(shifted, x);<br>gray.unroll(x, 4); |
| | Trial 4 | shifted.tile(x, y, xi, yi, 512, 32);<br>shifted.parallel(y);<br>gray.compute_at(shifted, x);<br>ly.compute_at(shifted, x);<br>gray.unroll(x, 4);<br>lx.compute_at(shifted, x); |
| Interp. | Initial Setup | // Initial schedule setup needed to compile<br>for (int l = 1; l < levels-1; ++l) {<br>   downsampled[l].compute_root();<br>   interpolated[l].compute_root();<br>} |
| | Trial 1 | downsampled[7].tile(x, y, xi, yi, 128, 8);<br>normalize.parallel(y);<br>downsampled[1].parallel(y);<br>interpolated[1].parallel(c);<br>interpolated[5].unroll(c, 3);<br>normalize.vectorize(x, 16);<br>downsampled[2].parallel(c);<br>interpolated[2].parallel(y);<br>downsampled[3].unroll(x, 2);<br>interpolated[1].vectorize(x, 4); |
| | Trial 2 | normalize.parallel(y);<br>interpolated[1].vectorize(x, 4);<br>downsampled[1].parallel(c);<br>downsampled[1].unroll(x, 4);<br>interpolated[1].parallel(c);<br>normalize.vectorize(x, 16);<br>downsampled[2].parallel(y);<br>downsampled[5].vectorize(x, 8);<br>interpolated[2].unroll(x, 4); |
| | Trial 3 | interpolated[2].tile(x, y, xi, yi, 256, 32); |

| | | normalize.parallel(y);<br>downsampled[1].parallel(y);<br>normalize.bound(x, 0, input.width());<br>interpolated[1].parallel(c);<br>interpolated[1].vectorize(x, 8);<br>normalize.unroll(x, 4);<br>downsampled[2].parallel(y);<br>normalize.unroll(x, 2);<br>downsampled[2].vectorize(x, 16);<br>interpolated[2].parallel(c);<br>downsampled[2].parallel(c); |
| | Trial 4 | downsampled[6].tile(x, y, xi, yi, 16, 8);<br>normalize.parallel(y);<br>downsampled[1].parallel(c);<br>normalize.vectorize(x, 16);<br>interpolated[1].parallel(c);<br>downsampled[2].tile(x, y, xi, yi, 64, 8);<br>downsampled[1].unroll(x, 4);<br>downsampled[2].parallel(y);<br>interpolated[2].parallel(y); |

**4. Best PPO Generated Schedule (GPU) for each Reinforcement Learning Trial:**

| | | |
|---|---|---|
| Blur | Trial 1 | blur_y.gpu_tile(x, y, xi, yi, 64, 8);<br>blur_y.unroll(yi, 4);<br>blur_y.unroll(xi, 2); |
| | Trial 2 | blur_y.gpu_tile(x, y, xi, yi, 128, 8);<br>blur_y.unroll(yi, 4);<br>blur_y.unroll(xi, 2); |
| | Trial 3 | blur_y.gpu_tile(x, y, xi, yi, 64, 8);<br>blur_y.unroll(xi, 2);<br>blur_y.unroll(yi, 4); |
| | Trial 4 | blur_y.gpu_tile(x, y, xi, yi, 64, 16);<br>blur_y.unroll(yi, 4);<br>blur_y.unroll(xi, 2); |
| Harris | Trial 1 | shifted.gpu_tile(x, y, xi, yi, 64, 16);<br>shifted.unroll(yi, 2);<br>shifted.unroll(yi, 4); |
| | Trial 2 | shifted.gpu_tile(x, y, xi, yi, 32, 16);<br>shifted.unroll(yi, 4);<br>shifted.unroll(yi, 2); |

| | | |
|---|---|---|
| | Trial 3 | shifted.compute_root();<br>shifted.gpu_tile(x, y, xi, yi, 8, 32);<br>shifted.bound(y, 0, input.height());<br>shifted.unroll(yi, 2);<br>gray.compute_at(shifted, x);<br>gray.gpu_threads(x, y);<br>gray.unroll(y, 2);<br>gray.unroll(x, 3); |
| | Trial 4 | shifted.gpu_tile(x, y, xi, yi, 32, 16);<br>shifted.unroll(yi, 4); |
| Interp. | Initial Setup | // Initial schedule setup needed to compile<br>for (int l = 1; l < levels-1; ++l) {<br>    downsampled[l].compute_root();<br>    interpolated[l].compute_root();<br>} |
| | Trial 1 | normalize.gpu_tile(x, y, c, xi, yi, ci, 8, 8, 3);<br>downsampled[1].gpu_tile(x, y, c, xi, yi, ci, 8, 8, 4);<br>interpolated[1].gpu_tile(x, y, c, xi, yi, ci, 16, 16, 4);<br>downsampled[2].gpu_tile(x, y, c, xi, yi, ci, 8, 16, 4);<br>downsampled[6].gpu_tile(x, y, c, xi, yi, ci, 16, 8, 4);<br>interpolated[3].gpu_tile(x, y, c, xi, yi, ci, 8, 8, 4);<br>downsampled[4].gpu_tile(x, y, c, xi, yi, ci, 16, 8, 4);<br>downx[5].compute_at(downsampled[5], x);<br>downsampled[8].gpu_tile(x, y, c, xi, yi, ci, 16, 16, 4);<br>interpolated[1].unroll(yi, 2);<br>interpolated[8].gpu_tile(x, y, c, xi, yi, ci, 16, 16, 4);<br>interpolated[2].gpu_tile(x, y, c, xi, yi, ci, 32, 8, 4);<br>interpolated[2].unroll(y, 3);<br>downsampled[5].gpu_tile(x, y, c, xi, yi, ci, 8, 32, 4);<br>downsampled[6].unroll(xi, 4);<br>downsampled[2].unroll(ci, 4);<br>downsampled[4].unroll(yi, 2);<br>interpolated[2].unroll(xi, 3);<br>downsampled[5].unroll(ci, 3);<br>downsampled[3].gpu_tile(x, y, c, xi, yi, ci, 8, 16, 4); |
| | Trial 2 | normalize.gpu_tile(x, y, c, xi, yi, ci, 32, 8, 3);<br>downsampled[1].gpu_tile(x, y, c, xi, yi, ci, 8, 16, 4);<br>interpolated[1].gpu_tile(x, y, c, xi, yi, ci, 8, 8, 4);<br>downsampled[2].gpu_tile(x, y, c, xi, yi, ci, 8, 8, 4);<br>interpolated[2].gpu_tile(x, y, c, xi, yi, ci, 16, 16, 4);<br>normalize.unroll(ci, 3);<br>downsampled[3].gpu_tile(x, y, c, xi, yi, ci, 8, 8, 4); |
| | Trial 3 | normalize.gpu_tile(x, y, c, xi, yi, ci, 8, 32, 3);<br>downsampled[1].gpu_tile(x, y, c, xi, yi, ci, 32, 8, 4);<br>interpolated[1].gpu_tile(x, y, c, xi, yi, ci, 16, 8, 4);<br>downsampled[2].gpu_tile(x, y, c, xi, yi, ci, 8, 8, 4); |

| | | |
|---|---|---|
| | | interpolated[2].gpu_tile(x, y, c, xi, yi, ci, 16, 16, 4);<br>downsampled[3].gpu_tile(x, y, c, xi, yi, ci, 8, 32, 4);<br>interpolated[5].unroll(x, 4);<br>interpolated[3].unroll(x, 4);<br>interpolated[6].unroll(y, 4);<br>normalize.unroll(ci, 2);<br>downsampled[1].unroll(xi, 2); |
| | Trial 4 | normalize.gpu_tile(x, y, c, xi, yi, ci, 8, 16, 3);<br>downsampled[1].gpu_tile(x, y, c, xi, yi, ci, 8, 8, 4);<br>interpolated[1].gpu_tile(x, y, c, xi, yi, ci, 16, 8, 4);<br>downsampled[2].gpu_tile(x, y, c, xi, yi, ci, 16, 8, 4);<br>downsampled[3].gpu_tile(x, y, c, xi, yi, ci, 8, 8, 4);<br>interpolated[6].unroll(c, 2);<br>interpolated[5].unroll(y, 3);<br>interpolated[2].gpu_tile(x, y, c, xi, yi, ci, 8, 16, 4);<br>interpolated[8].unroll(c, 3);<br>interpolated[3].gpu_tile(x, y, c, xi, yi, ci, 16, 16, 4);<br>downsampled[4].gpu_tile(x, y, c, xi, yi, ci, 16, 8, 4); |