

# AV1 Decoder Performance Optimization and Device Test

23-AUG-2022

**Presenter:** Ronald Bultje – Two Orioles Kaustubh Patankar - ITTIAM









# Introduction



### AV1 – Quick recap





# AV1 Decoder Optimization dAV1d



INLINE Dct8\_NEON(Residual \*buffer, step, bool transpose);

#### INVERSE TRANSFORMS:

GAV1 IMPLEMENTATION Gav1 has 2 versions of each inverse 1D transform with integrated load, transpose, store and residual add.

#### OVERHEAD

- Load/store overhead in 1D transforms
- Load/store overhead in transpose
- Load/store overhead in RowShift()
- Load overhead in residual\_add()
- Some overhead can be reduced by inlining, but that increases binary size

```
Dct8TransformLoopColumn_neon(Array2dView<Pixel> *frame,
Residual *buffer, tx_size, tx_type)
```



#### INVERSE TRANSFORMS:

GAV1 IMPLEMENTATION Gav1 has 2 versions of each inverse 1D transform with integrated load, transpose, store and residual\_add.

#### OVERHEAD

- Load/store overhead in 1D transforms
- Load/store overhead in transpose
- Load/store overhead in RowShift()
- Load overhead in residual add()
- Some overhead can be reduced by inlining, but that increases binary size

```
INLINE Dct8 NEON (Residual *buffer, step, bool transpose)
    [load]
    if (transpose) {
         [transpose]
      actual 1d inverse transform 1
    if (transpose) {
         [transpose]
    [store]
Dct8TransformLoopRow neon(Array2dView<Pixel> *frame,
                          Residual *buffer, tx size, tx type);
Dct8TransformLoopColumn neon(Array2dView<Pixel> *frame,
                             Residual *buffer, tx size, tx type);
```



#### INVERSE TRANSFORMS:

#### DAV1D IMPLEMENTATION

Dav1d uses transform implementations with a custom calling ABI which retains the vector registers as part of the function call interface. This is only possible because it's hand-written assembly (not intrinsics).

#### OVERHEAD

- Load/store overhead in 1D transforms
- Load/store overhead in transpose
- Load overhead in residual add()
- Load/store overhead in downshifts
- Custom calling ABI eliminates overhead without increasing binary size
- dav1d also eliminates the first transpose by integrating it with the scan table

:m_add_dct_dct_8x8_8bpc	_neon:
mov x15	5, x30
[ dc-only checks ]	
adr x5,	inv_dct_8h_x8_neon
adr x4,	inv_dct_8h_x8_neon
b inv	_txfm_add_8x8_neon
[ dc-only checks ] adr x5, adr x4, b <b>inv</b>	<pre>inv_dct_8h_x8_neo inv_dct_8h_x8_neo '_txfm_add_8x8_neon</pre>

#### inv\_txfm\_add\_8x8\_neon:

```
[ prologue ]
[ load coefficients ]
blr x4
[ 8x right-shift-by-one ]
[ transpose ]
blr x5
[ write out ]
[ epilogue ]
br x15
```

#### inv\_dct\_8h\_x8\_neon:

```
[ actual 1d transform ]
ret
```



#### **INVERSE TRANSFORMS:**

#### DAV1D IMPLEMENTATION

Dav1d uses transform implementations with a custom calling ABI which retains the vector registers as part of the function call interface. This is only possible because it's hand-written assembly (not intrinsics).

#### OVERHEAD

 Special identity implementation merges first 1D identity (which upshifts by 1) and downshifts (which downshift by 1)

#### inv\_txfm\_add\_identity\_dct\_8x8\_8bpc\_neon:

mov	x15, x30
adr	x5, inv_dct_8h_x8_neon
b	inv txfm identity add 8x8 neon

#### inv\_txfm\_identity\_add\_8x8\_neon:

```
[ prologue ]
[ load coefficients ]
[ transpose ]
blr x5
[ write out ]
[ epilogue ]
br x15
```

# inv\_dct\_8h\_x8\_neon: [ actual 1d transform ] ret





DCT/DCT Inverse Transforms

10



### other coding tool implementations

#### **CODING TOOLS:**

DAV1D vs. GAV1

- Loop Restoration: gav1's arm implementation of the selfguided filter is superior (because better cache efficiency)
- CDEF: both implementations have special-case implementations for single filters (when the primary or secondary CDEF filter strength is zero), but gav1's secondary-only filter implementation is (unexpectedly) slower than the dual-filter implementation
- Motion Vector Referencing SIMD: both decoders implement SIMD for different parts of the code
- Arithmetic coding: dav1d has hand-written assembly (including SIMD for CDF updates), which gav1 does not have
- Film Grain: both decoders have SIMD, but the Neon film grain is disabled in gav1
- Directional Intra Prediction (z1-3): gav1 has implementations for arm & x86, whereas dav1d only has SIMD implementations for x86

Conclusion: both decoders have some unique implementation ideas, and both can be improved further.



# Multi threading

#### MULTI-THREADING:

dav1d uses a task-queue design, where each component in the decoding loop runs as a generic task with a simple dependency resolution mechanism:

- entropy tile-sbrow reading
- tile-sbrow reconstruction
- deblock, cdef & loop restoration sbrow inloop post-filtering
- film grain sbrow out-loop post-filtering
- multiple frames in parallel

Together, this keeps low number of worker threads sufficiently busy without relying on particular bitstream features. Also, it does not require task-specific worker threads, and is therefore resource-friendly. Music Video (720p, 8-bit, 1000 frames) on arm (May, 2022)





# **Benchmarking and Analysis**



# Purpose of AV1 Benchmarking ON ANDROID ECOSYSTEM

- AV1 is the first generation of royalty-free video coding format developed by AOM
  - Helps to improve the quality of video delivery
- There are many service providers in the market already using AV1 for video streaming
- AV1 real-time playback capability on clients is critical
  - Already there is HW decoding support for AV1 on smart TV and set-top boxes
- Adoption of AV1 into android ecosystem is important and support for HW decoding is still not widely available and its penetration may take a couple of more years.
- Wide range of android devices with different chipsets are existing in the market and their processing capability varies
- This study helps to know the AV1 SW decoding capability on a sample of devices across the range and enables the market for wider adoption of AV1 delivery



SETUP AND CONFIGURATIONS





**MEASUREMENT METHODOLOGY** 

### Performance

- Using VLC Benchmarking APP
- Reports number of frames dropped
- Frame drop > 2% is non real time
- ST and MT

### **CPU** Load

- Using "adb top" command
- Validation of SW decoders
- Sampling at 1 sec interval



**DEVICE SELECTION** 

Category	Cortex Version	ARMArchitecture	Clock Frequency	
High-end	4xA7x + 4xA5x	v8	Any	
	2xA7x + 6xA55	v8	Any	
Mid-range	8xA53	v8	>= 1.8GHz	
	2xA7x + 6xA53	v8	Any	
Low-end	8xA53	v8	< 1.8 GHz	
	Any	٧7	Any	

- Total devices Android 61
- Covering 15 global brands
- Release year from 2018 to 2020
- Classification as High, Medium and

Low complexity

• Based on CPU capability



TEST CLIP, CONFIGURATIONS,

- TOOLS ON [T-ON]: all AV1 supported coding tools are used while encoding.
- TOOLS OFF[8] [T-OFF]: AV1 coding tools with a poor trade-off between decoding complexity and compression efficiency are disabled intelligently (fast-decode option in SVT-AV1) as described below.
  - Enable Motion Field MV prediction only for a well-predictable motion field
  - Enabled tools
    - o OBMC
    - $\circ$  Warp
    - $\circ$  Partition depth
    - Reduced complexity approximations on Self-Guided filter, CDEF
  - Intelligent application of CDEF and Deblocking at super block level

### Test Clip – AOM – CTC

- Original resolution of 1920x1080
- Down scaled to 5 resolutions 270p, 360p, 540p, 720p, 1080p
- FPS-30 and 60

• QP Points –	9
---------------	---

18

QP	Average Bitrate					
	Tools-On			Tools-Off		
	54op	720p	108op	54op	720p	108op
32	2.08	2.99	5.12	2.06	3.02	5.14
43	1.05	1.46	2.47	1.04	1.46	2.46
55	0.45	0.61	1.06	0.45	0.61	1.04
6 <sub>3</sub>	0.12	0.17	0.3	0.12	0.16	0.3



### Standalone Decoder Performance







### **Tools ON vs Tools OFF**

### DAV1D WITH SINGLE AND 4 THREADS



**QP** | Resolution



### Performance of David and Gavi

### SINGLE AND 4 THREADS – TOOLS OFF CONFIGURATION



QP | Resolution



### Average CPU Utilization

### DAV1D AND GAV1 – SINGLE AND FOURTHREADS





### Conclusion





