



# PowerVR

## Performance Recommendations

Copyright © Imagination Technologies Limited. All Rights Reserved.

This publication contains proprietary information which is subject to change without notice and is supplied 'as is' without warranty of any kind. Imagination Technologies and the Imagination Technologies logo are trademarks or registered trademarks of Imagination Technologies Limited. All other logos, products, trademarks and registered trademarks are the property of their respective owners.

Filename : PowerVR.Performance Recommendations  
Version : PowerVR SDK REL\_3.5@3523383a External Issue  
Issue Date : 17 Apr 2015  
Author : Imagination Technologies Limited

# Contents

<b>1. Introduction .....</b>	<b>4</b>
1.1. Document Overview .....	4
1.2. The Golden Rules .....	4
1.3. Optimal Development Approach .....	4
1.4. Understanding Rendering Bottlenecks .....	4
<b>2. Optimizing Geometry .....</b>	<b>6</b>
2.1. Geometry Complexity .....	6
2.2. Primitive Type .....	6
2.3. Data Types .....	6
2.3.1. "Fixed" Data Types .....	6
2.4. Interleaving Attributes .....	6
2.5. Vertex Buffer Objects .....	7
2.6. Padding .....	7
<b>3. Optimizing Textures .....</b>	<b>8</b>
3.1. Texture Size .....	8
3.2. Texture Compression .....	8
3.2.1. Why use PVRTC? .....	8
3.2.2. Image File Compression vs. Texture Compression .....	9
3.3. MIP-Mapping .....	11
3.3.1. Advantages .....	11
3.3.2. Generation .....	11
3.3.3. Filtering .....	11
3.4. Texture Sampling .....	12
3.4.1. Texture Filtering .....	12
3.4.2. Dependent Texture Reads (Series5 and Series5XT only) .....	12
3.4.3. Wide Floating Point Textures .....	12
3.5. Demystifying NPOT .....	13
3.5.1. PowerVR Series5 Support .....	13
3.5.2. GL_IMG_texture_npot .....	13
3.5.3. Guidelines .....	13
3.6. Texture Uploading .....	14
3.6.1. Texture Warm-up .....	14
3.6.2. Texture Formats and Precision .....	14
3.7. Render to Texture .....	14
3.8. Mathematical Look-ups (Series5 and Series5XT only) .....	14
<b>4. Optimizing Shaders .....</b>	<b>15</b>
4.1. Choose the Right Algorithm .....	15
4.2. Know Your Spaces .....	15
4.3. Flow Control .....	15
4.4. Demystifying Precision .....	16
4.4.1. Highp .....	16
4.4.2. Medump .....	16
4.4.3. Lowp (Series5 and Series5XT only) .....	16
4.4.4. Attributes .....	16
4.4.5. Varyings .....	17
4.4.6. Samplers .....	17
4.4.7. Uniforms .....	17
4.4.8. Conversion Costs .....	17
4.5. Scalar Operations .....	18
4.6. Sparse Matrices .....	18
4.7. "Const" Data in Shaders .....	18
<b>5. Optimizing Specific Techniques .....</b>	<b>19</b>
5.1. Multiple Render Targets (Series6 only) .....	19
5.2. Efficient Sprite Rendering .....	19

6.	Contact Details .....	20
----	-----------------------	----

## List of Figures

Figure 1. Cyclical profiling .....	4
Figure 2. Image file compression vs. texture compression.....	10
Figure 3. Increasing complexity and reducing processing .....	19

# 1. Introduction

PowerVR Series5 and Series6 are families of Graphics Cores from Imagination Technologies designed specifically for shader-based APIs like OpenGL ES 2.0 and 3.0. Due to their scalable architectures, the PowerVR family spans a huge performance range.

## 1.1. Document Overview

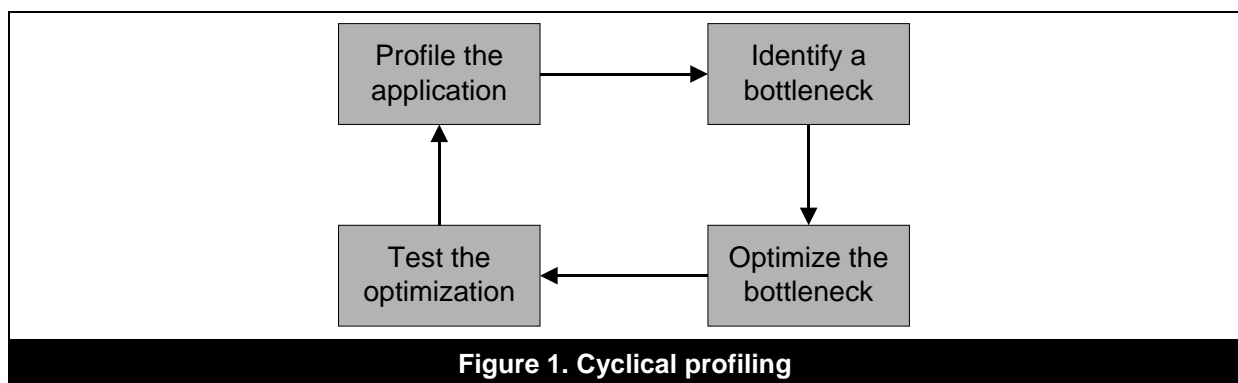
The purpose of this document is to serve as recommendation and advice for developers who wish to get the best graphics performance from a PowerVR Series5 or Series6 enabled device. Throughout the document, the specific recommendations for PowerVR Series5 and Series6 are marked as appropriate.

## 1.2. The Golden Rules

The golden rules are a set of more generic performance recommendations that developers should seek to implement and observe as many of the techniques and principles mentioned in these rules help to produce well-behaved, high performance graphics applications. These rules are detailed in the document entitled “PowerVR Performance Recommendations: The Golden Rules”, which is supplied with the PowerVR SDK.

## 1.3. Optimal Development Approach

It is crucial to adopt the practices identified in this document from the very start of development in order to save much time and effort later. Once an application is implemented to a near-final state, the process of iteration depicted in Figure 1 should be adopted. The main benefit of this approach is that time is not wasted and graphics quality is not comprised by making changes that do not benefit performance.



## 1.4. Understanding Rendering Bottlenecks

It is a common misconception that the same actions can speed up any application. For example:

- **Polygon count reduction:** If the bottleneck of the application is fragment processing or texture bandwidth then the only result of this action will be to reduce the graphical quality of the application without improving rendering speed. In fact, if simpler models cause more of the render target to be covered by a material with complex fragments then this can actually slow down an application.
- **Reduce rendering resolution:** In this case, if the fragment processing workload of your application is not the bottleneck then this will also only serve to reduce the quality of the graphics in your application without improving performance.

In reality, it is only once the limiting factor of an application is determined by profiling with the correct tools that optimization work should be applied. It is also important to realise that once work has been

done then the application requires re-profiling in order to determine whether the work actually improved performance and whether the bottleneck is still at the same stage of the graphics pipeline. It may be that the limiting stage in rendering is now at a different place and further optimization should be targeted accordingly.

## 2. Optimizing Geometry

### 2.1. Geometry Complexity

It is important that an appropriate level of geometry complexity be used for each object or portion of an object. It is a waste to use a large number of polygons on an object that will never cover more than a small area of the screen. Likewise, it is a waste to use polygons for detail that will never be seen due to camera angle, or culling, or to use large amounts for objects that may be drawn with much fewer (such as spending hundreds of polygons drawing a single quad). Shader techniques such as bump mapping should be considered to minimize geometry complexity, but still maintain a high level of perceived detail. Techniques such as “Level of Detail” should be used. This is especially true for things such as reflection passes where higher amounts of geometry may not be visible.

### 2.2. Primitive Type

In general, drawing a mesh as a single, indexed, triangle list will ensure the best performance on PowerVR Series5 hardware.

### 2.3. Data Types

Vertex shaders always expect attributes to be of the type `float`, this means that all types except `float` will require a conversion. This conversion is performed in the USSE pipeline and costs a few additional cycles. Thus the choice of attribute data type is a trade-off between shader cycles, bandwidth/storage requirements and precision. It is important that type conversion is considered as bandwidth is always at a premium.

Precision requirements should be checked carefully, the `byte` and `short` types are often sufficient, even for position information. For example, scaled to a range of 10m the `short` types give a precision of 150  $\mu\text{m}$ . Scaling and biasing those attribute values to fit a certain range can often be folded into other vertex shader calculations, e.g., multiplied into a transformation matrix.

#### 2.3.1. “Fixed” Data Types

The `fixed` data type uses the same bandwidth as `float`, but requires additional format conversion cycles in the USSE pipeline, thus it should be avoided.

### 2.4. Interleaving Attributes

Two ways exist to store vertex data in memory, either the data is stored with all the information, position, normals, etc., pertaining to a given vertex in a single block, followed by all the information pertaining to the next vertex, and so on, or the data can be stored in a series of arrays, each containing all the information of a particular type for each vertex. For example, an array of positions, an array of normals, etc. The first of these two options is called “interleaving”. In general data should be interleaved as this provides better cache efficiency, and thus better performance.

Two major caveats exist to this rule. Interleaving should not be used if several meshes are to share the same array of vertex attributes. In this case putting the instances of this attribute into their own array may result in better performance, and will save bandwidth and storage space due to there being less duplication.

Interleaving should also not be used if a single attribute will be updated frequently, outside of the Graphics Core, while the other attributes remain the same. In this instance, data that will not be updated should be interleaved, while data that will be updated is held in a separate array.

## 2.5. Vertex Buffer Objects

Vertex Buffer Objects (VBO) (and, where available, Vertex Array Objects) are the preferred way of storing vertex and index data. Since VBO storage is managed by the driver there is no need to copy an array from the client side at every draw call and the driver is able to perform some transparent optimizations.

Pack all the vertex attributes that are required for a mesh into the same VBO unless a mixture of static and dynamic attributes are being used. Do not create a VBO for every mesh, it is a good idea to group meshes that are always rendered together in order to minimize buffer rebinding, this also has the benefit of improving batching.

As the TBDR tends to process multiple frames at a time, the driver has to internally allocate multiple buffers for dynamic VBOs so that each frame has a unique dynamic buffer associated with it. Because dynamic VBOs cause the driver to behave in this way it is generally better for performance to resubmit vertex data that changes on a per-frame basis. If there is a mesh where only some of the vertex data is dynamic (for example, a skinned character in a game) then a VBO should be created that contains the static data and use calls to `glVertexAttribPointer()` to resubmit the dynamic vertex data. On a similar note, a VBO that will never change should always set `STATIC_DRAW` while a VBO whose contents will change should never set it.

## 2.6. Padding

When vertex data is interleaved, each vertex should be aligned to a four byte boundary. When vertex data is not interleaved each element in each array of vertex data should be aligned to a four byte boundary.

## 3. Optimizing Textures

### 3.1. Texture Size

It is a common misconception that bigger textures always look better; a 1024x1024 texture that never takes up more than a 32x32 area of the screen is a waste of both storage space and reduces cache efficiency. A texture's size should be based on its usage; there should be a 1 pixel to 1 texel mapping when the object that it is mapped to is viewed from the closest allowable distance.

In some instances it is acceptable to use images that are bigger than the number of texels the texture will cover on the screen. This is primarily in situations where low bit rate texture compression is being used, but the quality of the compressed texture is not deemed to be high enough. In these instances using a larger compressed texture may produce a more acceptable image quality, while still ensuring bandwidth savings over uncompressed textures.

### 3.2. Texture Compression

Modern applications have become graphically intensive. Certain types of software, such as games or navigation aids, often need large amounts of textures in order to represent a scene with satisfying quality. Texture compression can save or allow better utilization of bandwidth, power, and memory without noticeably losing graphical quality and should be used as much as possible. PowerVR hardware offers a specific form of texture compression called "PVRTC" which should be used as much as possible.

PVRTC is PowerVR's proprietary texture compression scheme. It uses a sophisticated amplitude modulation scheme to compress textures: texture data is encoded as two low-resolution images along with a full resolution, low bit-precision modulation signal. More information can be found in the whitepaper:

Fenney, S. (2003) 'Texture Compression Using Low-Frequency Signal Modulation' *SIGGRAPH Conference*.

Additionally, it supports both opaque (RGB) and translucent (RGBA) textures (unlike other formats, such as S3TC, that require a dedicated, larger form to support full alpha channels). It also boasts a very high image quality for competitive compression ratios: 4 bits per pixel (PVRTC 4bpp) and 2 bits per pixel (PVRTC 2bpp). At time of writing, no other format is available in hardware at such a low bit rate.

#### 3.2.1. Why use PVRTC?

In any given situation, the best texture format to use is the one that gives the required image quality at the highest rate of compression. The smaller the size of the texture data, the less bandwidth is required for texture fetches; this reduces power consumption, can increase performance, and allows for more textures to be used for the same budget. The smallest RGB and RGBA format currently available is PVRTC 2bpp and, as such, it should be considered for every texture in an application. Larger formats (such as PVRTC 4bpp) should only be used if the image quality provided by a particular PVRTC 2bpp image does not have sufficient quality.

#### Storage Footprint vs. Memory Footprint

PVRTC compression reduces the memory footprint of a given texture. This allows applications to fit all their required textures in a constrained amount of texture memory, or to use larger (or more) textures for the same memory budget at, potentially, extra quality. In addition, any savings in memory requirements are very useful for mobile and tablet devices where memory is shared across an entire SoC (System on Chip).

#### Performance Improvement

The smaller memory footprint of PVRTC means less data is transferred from memory to the Graphics Core allowing for major bandwidth savings. In situations where memory bandwidth is the limiting factor in an application's performance PVRTC can provide a significant boost.



## Power Consumption

Memory accesses are one of the primary causes of increased power consumption on mobile devices where battery life is of the utmost importance. The bandwidth savings and better cache performance resulting from the use of PVRTC both contribute to decreasing the quantity and magnitude of memory accesses; which in turn reduce the power consumption of an application.

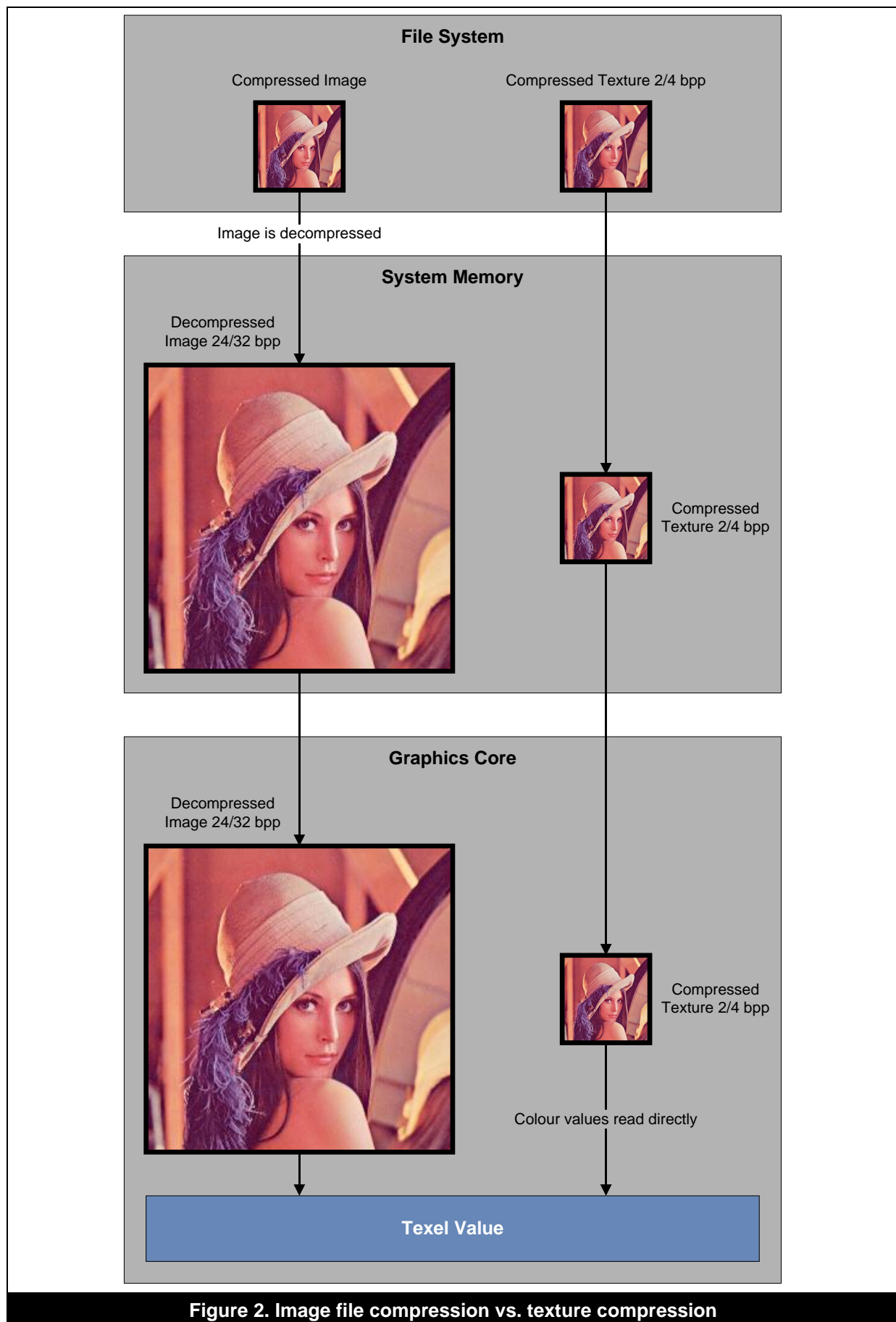
### 3.2.2. Image File Compression vs. Texture Compression

Developers are familiar with compressed image file formats such as JPG or PNG. It is important to be aware of the distinction between these forms of “storage” compression and the texture compression discussed in this document.

The primary requirement of storage compression schemes is that files compressed using them should occupy as small an amount of storage in a file system as possible. There is no requirement that the data stay compressed while in use. The result is that storage-based image file formats tend to produce very small file sizes, often for very high (or lossless) image quality, but at the cost of immediate decompression on use. This immediate decompression, usually to 24/32bpp means that the image, while small on disk, consumes large amounts of bandwidth and memory at runtime.

Texture compression schemes, such as PVRTC are designed to be directly usable by the Graphics Core. The texture data exists in storage, in memory, and when transferred to the graphics hardware itself, in the compressed format. The only step in which full-precision colour values are extracted from a compressed state is when dedicated texture sampling hardware inside the graphics accelerator passes texel values to the shader processing units. A graphical representation of this can be seen in Figure 2.

This allows all the advantages mentioned in Section 3.2.1, but puts some limits on the form the compression technique may take. In order to allow for direct use by the graphics accelerator a texture format should be optimized for random access, with a minimal size of data from which to retrieve each texel’s values. Consequently, texture compression schemes are usually fixed bitrate with very high data locality. Image file formats are not constrained by these requirements and thus can often achieve higher compression ratios and image quality for a given data size.



### 3.3. MIP-Mapping

MIP-maps are smaller, pre-filtered variants of a texture image, representing different levels-of-detail of a texture. By using a minification filter mode that uses MIP-maps, the Graphics Core can be set up to automatically calculate which level-of-detail comes closest to mapping the texels of a MIP-map to pixels in the render target, and use the right MIP-map for texturing.

#### 3.3.1. Advantages

Using MIP-maps has two important advantages, namely it increases performance by massively improving texture cache efficiency, especially in cases of strong minification. It also improves image quality by countering the aliasing that is caused by the under-filtering of textures when MIP-mapping. The single drawback of MIP-mapping is that it requires approximately 1/3 more texture memory per image. Depending on the situation, this cost may be minor when compared to the benefits in terms of rendering speed and image quality.

There are some exceptions where MIP-maps should not be used. Specifically, MIP-mapping should not be used where filtering cannot be applied sensibly, such as for textures that contain non-image data such as indices or depth textures. It should also be avoided for textures that are never minified, for example, UI elements where texels are always mapped one-to-one to pixels.

#### 3.3.2. Generation

Ideally MIP-maps should be created offline using a tool like PVRTexTool (available as part of the PowerVR Graphics SDK). It is, however, possible to generate MIP-maps at runtime using the function `glGenerateMipmap` and this can be useful for updating the MIP-maps for a render to texture target. This will not work, however, with PVRTC textures which must have their MIP-maps generated offline. A decision must be made as to which cost is the most appropriate, the storage cost of offline generation, or the runtime cost of `glGenerateMipmap`.

#### 3.3.3. Filtering

Finally, it should be noted that the lack of filtering between MIP-map levels can lead to visible seams at MIP-map transitions, a form of artifacting called “MIP-map banding”. Trilinear filtering using the filter mode `GL_LINEAR_MIPMAP_LINEAR` can effectively eliminate these seams, for a price (see Section 3.4.1), and thus achieve an even higher image quality.

## 3.4. Texture Sampling

### 3.4.1. Texture Filtering

Texture filtering is used to increase the image quality of textures used in 3D scenes. However, it comes at a cost. Filtering works by taking multiple texture fetch values and combining them in order to produce as good a sampling value as possible to use in fragment calculations. Retrieving multiple values requires more data to be fetched, possibly from disparate areas of memory and so cache performance and bandwidth use can be affected. For instance, whenever two MIP-map levels must be blended together for trilinear filtering, the texture unit in the Graphics Core must spend time and bandwidth fetching and filtering the required data from the two MIP-map levels in question. This can cause the processing of a fragment to stall while the data is fetched and adds to the total amount of memory that must be transferred across the bus in order to render a frame.

For independent texture reads, texture sampling can begin before the execution of a shader and so the latency of the texture fetch can be avoided. For dependent reads the cost can further be amortised thanks to the hardware scheduler in PowerVR Graphics Cores, particularly if the shader in question involves a lot of mathematical calculation. This latency can be hidden by swapping in another thread on the Graphics Core. This thread will process as much as possible with the original thread being swapped back once the fetch is complete. Further information on the functioning of the Coarse Grain Scheduler and thread scheduling within PowerVR hardware can be found in the “PowerVR Hardware Architecture Guide for Developers”.

The three main techniques for texture filtering are bilinear, trilinear, and anisotropic, where each gives increased image quality than the previous, at an increasing cost. Performance can be gained by using an appropriate level of filtering, following the principle of “good enough” (see “PowerVR Performance Recommendations: The Golden Rules”). Also, not using anisotropic if trilinear is acceptable. Not using trilinear if bilinear is also acceptable. It is recommended not to filter at all if it is not necessary (choose nearest or point texture sampling).

### 3.4.2. Dependent Texture Reads (Series5 and Series5XT only)

A dependent texture read is a texture read in which the texture coordinates depend on some calculation within the shader instead of on a varying. As the values of this calculation cannot be known ahead of time it is not possible to pre-fetch texture data and so stalls in shader processing occur.

Vertex shader texture lookups always count as dependent texture reads, as do texture reads in fragment shaders where the texture read is based on the `.zw` channels of a varying. On some driver and platform revisions `Texture2DProj()` also qualifies as a dependent texture read if given a `Vec3` or a `Vec4` with an invalid `w`.

The cost associated with a dependent texture read can be amortised to some extent by hardware thread scheduling, but they should still be avoided wherever possible for good performance. On PowerVR Series6 Graphics Cores this is no longer a concern as there is no additional cost to a dependent texture read.

### 3.4.3. Wide Floating Point Textures

For textures that exceed 32 bits per texel, each additional 32 bits is counted as a separate texture read. This also applies to half float texture with 3 or 4 components as well as float textures with 2 or more components. These larger formats should be avoided unless necessary for a particular effect.

## 3.5. Demystifying NPOT

If a 2D texture has dimensions which are a power-of-two (i.e., width and height are  $2^n$  and  $2^m$  for some  $m$  and  $n$ ), then the texture is said to be a POT texture (power-of-two). If they are not it is said to be an NPOT texture (non-power-of-two). This section seeks to clarify the status of NPOT textures on PowerVR Series5 cores.

### 3.5.1. PowerVR Series5 Support

NPOT textures are supported as required by the OpenGL ES specifications. However, it is necessary to point out the following:

- NPOT textures are not supported in OpenGL ES 1.1 implementations.
- NPOT textures are supported in OpenGL ES 2.0 implementations, but only with the wrap mode of `GL_CLAMP_TO_EDGE`.
  - The default wrap mode in OpenGL ES 2.0 is `GL_REPEAT`. This must be specifically overridden in an application to `GL_CLAMP_TO_EDGE` for NPOT textures to function correctly.
  - If this wrap mode is not correctly set then an “invalid texture” error will occur, likewise a driver error may occur at runtime, on newer drivers, to highlight the need to set a wrap mode.

### 3.5.2. GL\_IMG\_texture\_npot

An extension exists (`GL_IMG_texture_npot`) to provide some of the functionality found outside of the core OpenGL ES specification. This extension allows the use of the following filters for NPOT textures:

- `LINEAR_MIPMAP_NEAREST`
- `LINEAR_MIPMAP_LINEAR`
- `NEAREST_MIPMAP_NEAREST`
- `NEAREST_MIPMAP_LINEAR`

It also allows the calling of `glGenerateMipmapOES` with an NPOT texture to generate NPOT MIP-maps. Like all other OpenGL extensions, the application should check for this extension’s presence before attempting to load and use it.

### 3.5.3. Guidelines

Finally, a few additional points should be considered when using NPOT textures:

- POT textures should be favoured over NPOT textures for the majority of use cases as this gives the best opportunity for the hardware and driver to work optimally.
- A 512x128 texture will qualify as a POT texture, not an NPOT texture, where rectangular POT textures are fully supported.
- 2D applications (such as a browser or other application rendering UI elements where an NPOT texture is displayed with a one-to-one texel to pixel mapping) should see little performance loss from the use of NPOT textures other than possibly at upload time.
- To ensure that texture upload can be optimally performed by the hardware, use textures where both dimensions are multiples of 32 pixels.
- The use of NPOT textures may cause a drop in performance during 3D rendering. This can vary depending upon MIP-map levels, size of the texture, texture usage and the target platform.

## 3.6. Texture Uploading

When a texture is uploaded through the use of `glTexImage2D` the input data is usually in linear scan-line format. Internally, PowerVR hardware uses its own layout to improve memory access locality and improve cache efficiency. Reformatting of the data is done on chip by dedicated hardware and thus is very fast, however, it is still recommended that a few steps be taken to minimize the cost of this reformat.

- Textures should be uploaded during non-performance critical periods, such as initialisation. This helps avoid the frame rate dips associated with additional texture loading.
- Avoid uploading texture data mid-frame to a texture object that has already been used for that frame.
- Consider performing a “warm-up” step after texture uploads have been performed. Once again, this helps avoid the frame rate dips associated with texture loading.

### 3.6.1. Texture Warm-up

The warm-up step mentioned before ensures that textures are fully uploaded immediately. By default, `glTexImage2D` does not perform all the processing required to upload immediately. Instead, the texture is fully uploaded the first time it is used. It is possible to force an upload by drawing a series of triangles off screen or otherwise obscured with the texture object in question bound and so marked for use. Performing this for all textures in a scene will avoid the cost and potential stutters when they are uploaded on first use.

### 3.6.2. Texture Formats and Precision

In general, textures should be read at `lowp` (see Section 4.4.6). The exceptions to this are half float textures which should be read as `mediump`, and float and depth textures which should be read as `highp`.

## 3.7. Render to Texture

The preferred method for rendering to textures on OpenGL ES 2.0 is through the use of Frame Buffer Objects (FBOs) with textures as attachments. The only situation where FBOs are not recommended is when accessing render targets from the CPU.

For maximised performance, FBOs should be rendered to in series, submitting all calls for one FBO before moving to the next. This serves to minimize state changes, as well as reducing unnecessary memory bandwidth usage caused by flushing partially completed renders when the target FBO is changed. Reuse of FBO targets should be avoided, where, if the previous content of an FBO is required for some stage of a render that has not been completed then this content will need to be preserved in an expensive copy operation. Due to this, a technique that may seem to be saving memory could actually be using substantially more. In addition, this reuse may serialise the processing of the two renders, which will cause a drop in performance.

## 3.8. Mathematical Look-ups (Series5 and Series5XT only)

Sometimes it can be a good idea to encode the results of a complex function into a texture and use it as a look-up table instead of performing the calculations in a shader. However, this will only provide a boost in performance if a bottleneck has been identified in the processing of the shader in question, and bandwidth is free to perform the texture lookup. If the function parameters (and thus the texture coordinates in the look-up table) vary wildly between adjacent fragments then cache efficiency will suffer. This can be mitigated through the use of MIP-maps but at the cost of accuracy. As a significant amount of work must be saved for this to be an optimisation, profiling should be performed to determine if the results of using look-up tables are acceptable. Due to the higher ALU power in Series6 onwards it is unlikely this will provide a benefit.

## 4. Optimizing Shaders

### 4.1. Choose the Right Algorithm

For complex shaders that run for more than a few cycles, picking the right algorithm is usually more important than low-level optimizations. It is highly recommended that a fast, well designed, algorithm be favoured over small performance tweaks to a poor algorithm. Bear in mind, that, although increasingly powerful, mobile graphics hardware is not designed to handle some of the latest techniques in desktop and console shaders. As such, a reduction in complexity will likely be needed from some of these techniques for mobile shader implementations.

### 4.2. Know Your Spaces

A common mistake in vertex shaders is to perform unnecessary transformations between model space, world space, view space and clip space. If the model-world transformation is a rigid body transformation, i.e., it only consists of rotations, translations, and mirroring, lighting and similar calculations can be performed directly in model space. Transforming uniforms such as light positions and directions to model space is a per-mesh operation, as opposed to transforming the vertex position to world or view space once per vertex and so is an optimization. In cases where a particular space must be used, e.g., for cube map reflections, it is often best to use this single space throughout.

### 4.3. Flow Control

PowerVR hardware offers full support for flow control in both vertex and fragment shaders without the need to explicitly enable an extension. When conditional execution depends on the value of a uniform variable, this is called “static flow control”, and the same shader execution path is applied to all vertex or fragment instances in a draw call. “Dynamic flow control”, on the other hand, refers to conditional execution based on per-fragment or per-vertex data, e.g., textures or vertex attributes.

Static flow control can be used to combine many shaders into one big “uber-shader”. Thorough profiling should be done when taking this approach, however, as a performance advantage may not be gained. A better solution when an uber-shader is desired is to use pre-processor defines to create separate shaders from one larger shader at build time, effectively creating many smaller shaders from a single original source file.

Using dynamic branching in a shader has a non-constant overhead that depends on the exact shader code. Dynamic branching is, therefore, unpredictable in its effect on performance. In general, the following specific points should be considered:

- Make use of conditionals to skip unnecessary operations when the condition is met in a significant number of cases.
- If the product of two complex functions is required, and that product sometimes evaluates to zero, use the less complex function, or the function that most often returns zero, as a condition for executing the other.
- Do not branch to `discard` (see “PowerVR Performance Recommendations: The Golden Rules”).
- **Series5 and Series5XT only:** Avoid branching to a texture read as samplers in dynamic branches qualify as “dependent texture reads” and will harm performance.
- **Series5 and Series5XT only:** The branching granularity on PowerVR Series5 hardware is one fragment or one vertex, meaning that area of fragments do not have to be spatially coherent in terms of branching.
- **Series6 only:** Unlike the previous PowerVR Series5 hardware, for branching to be effective on PowerVR Series6, fragments that branch should be spatially coherent.



## 4.4. Demystifying Precision

PowerVR hardware is designed with support for the multiple precision features of graphics APIs such as OpenGL ES 2.0 and OpenGL ES 3.0. Three precision modifiers are included in the API spec for OpenGL ES 2.0 onwards, namely `mediump`, `highp`, and `lowp`. Lower precision calculations can be performed faster, but need to be used carefully to avoid trouble with visible artefacts being introduced. The best method of arriving at the right precision for a given value is to begin with `lowp` or `mediump` for everything (except samplers) then increase the precision of specific variables until the visual output is as desired.

### 4.4.1. Highp

Float variables with the `highp` precision modifier will be represented as 32 bit floating point values, whereas integer values range from  $2^{31}-1$  to  $-2^{31}$ . This precision should be used for all vertex position calculations, including world, view, and projection matrices, as well as any bone matrices used for skinning where the precision, or range, of `mediump` is not sufficient. It should also be used for any scalar calculations that use complex built-in functions such as `sin`, `cos`, `pow`, `log`, etc.

### 4.4.2. Mediump

Variables declared with the `mediump` modifier are represented as 16 bit floating point values covering the range  $[65520, -65520]$ . The integer values cover the range  $[2^{15}-1, -2^{15}]$ . This precision level typically offers a performance improvement over `highp`, and should be considered wherever `highp` would normally be used (provided the precision is sufficient and maximum and minimum values will not be overflowed).

### 4.4.3. Lowp (Series5 and Series5XT only)

A variable declared with the `lowp` modifier will use a 10 bit fixed point format, allowing values in the range  $[-2, 2]$  to be represented to a precision of  $1/256$ . The integer values are in the range of  $[2^9-1, -2^9]$ . This precision is useful for representing colours and any data read from low precision textures, such as normals from a normal map. Care must be taken not to overflow the maximum or minimum value of `lowp` precision, especially with intermediate results.

#### Swizzling

Swizzling is the act of accessing or reordering the components of a vector out of order. Some examples of swizzling can be found next:

```
a = var.brg;           // Swizzled - Out of order access
b = vec3(var.g, var.b, var.r); // Swizzled - Out of order access
c = vec3(vec4);        // Not swizzled - Dropping a component does not change
                        // access order
d.gr = a.gr + b.gr     // Not swizzled - This will be optimized to a
                        // non-swizzled form
```

Swizzling costs performance when performed on `lowp` variables due to the additional work required to move vector components when they in `lowp` form, and thus should be avoided. PowerVR Series6 works exclusively on 16bit floating point and 32bit floating point scalar values, as such `lowp` will always be honoured as `mediump` and swizzling will no longer have any performance impact.

### 4.4.4. Attributes

The per-vertex attributes passed to a vertex shader should use a precision appropriate to the data-type being passed in, so, for example, `highp` would be unrequired for a float whose maximum value never goes above 2 and for which a precision of  $1/256$  would be acceptable.



#### 4.4.5. Varyings

Varyings represent the outputs from the vertex shader which are interpolated across a triangle and then fed into the fragment shader. Each varying requires additional space in the parameter buffer, and additional processing time to perform interpolation. To keep this to a minimum, as few a number of varyings as possible should be used.

##### Packing Varyings

Packing multiple varyings together, for example packing two `Vec2` into a single `Vec4` should suffer no performance penalty and will save varyings. Exclusively on PowerVR Series5 and Series5XT, coordinate varyings which are packed into the `.zw` channel of a `Vec4` will always be treated as a dependent texture read and should be avoided (see Section 3.4.2).

#### 4.4.6. Samplers

Samplers are used to sample from a texture bound to a certain texture unit. The default precision for sampler variables is `lowp`, and generally this is good enough. Two main exceptions exist to the `lowp` rule. If the sampler will be used to read from either a depth or float texture then it should be declared with `highp`. On the other hand, if the sampler will be used to read from a half float texture then it should be declared as `mediump`.

#### 4.4.7. Uniforms

Uniform variables represent values that are constant for all vertices or fragments processed as part of a draw call. Similar to redundant state changes, redundant uniform updates in between draw calls should be avoided. Unlike attributes and varyings, uniform variables can be declared as arrays. However, care should be taken when using uniform arrays. This is because while a certain number of uniforms can be stored in registers on-chip, large uniform arrays will be stored in memory and accessing them comes at a bandwidth and execution time cost.

##### Uniform Calculations

The PowerVR shader compiler is able to extract calculations based on uniforms from the shader and perform these calculations once per draw call. If this functionality is desired, it is important that the order of operations is chosen so that the uniforms are processed first, such as in the next example.

```
uniform highp mat4 modelview, projection;
attribute vec4 modelPosition;

// Can be extracted
gl_Position = (projection * modelview) * modelPosition;

// Cannot be extracted
gl_Position = projection * (modelview * modelPosition);
```

#### 4.4.8. Conversion Costs

When performing arithmetic on multiple precisions within the same calculation it is likely that values will have to be “packed” or “unpacked”. Packing is the act of taking a higher precision value and placing into a lower precision variable while unpacking is the reverse and involves taking a lower precision value and placing it into a higher precision variable.

Where possible precisions should be kept the same for an entire calculation as each pack and unpack has a cost associated with it. This cost can be further amortised by writing shaders in such a way that all higher precision calculations are performed together, at the top of the shader, and all lower precision calculations performed at the bottom. This ensures that variables are not repeatedly packed and unpacked. It also ensures that variables are not all unpacked into `highp` thereby losing any benefit of using lower precision.

## 4.5. Scalar Operations

It is very easy to accidentally vectorise a calculation. Hence, one should be wary of vectorising scalar operations where it cost more cycles for the same output. For example:

```
highp vec4 v1, v2;
highp float x, y;

// Bad
v2 = (v1 * x) * y; // vector * scalar followed by vector * scalar totals 8 scalar muladds

// Good
v2 = v1 * (x * y); // scalar * scalar followed by vector * scalar totals 5 scalar muladds
```

## 4.6. Sparse Matrices

If it is already known that many elements of a transformation matrix are zero, do not perform a full matrix transform. For example, given a typical projection matrix in the form of:

A	0	0	0
0	B	0	0
0	0	C	D
0	0	E	0

If the vertex being transformed is already in view space, an additional full transformation would be both a waste of cycles and unnecessary since dividing the matrix by a positive constant will not change the transformation result in homogeneous coordinates. In this case it is sufficient to store just four values. Similarly, non-projective transformation matrices usually have the fourth row fixed at (0, 0, 0, 1). If this holds true then it is possible to store the matrix as three `vec4` rows, replacing the matrix-vector multiplication with three dot products as shown next:

```
attribute highp vec3 vertexPos;
uniform highp vec4 modelview[3]; // first three rows of modelview matrix
uniform highp vec4 projection;   // = vec4(A/D, B/D, C/D, E/D)

void main()
{
    // transform from model space to view space
    highp vec3 viewSpacePos;
    viewSpacePos.x = dot(modelview[0], vec4(vertexPos, 1.));
    viewSpacePos.y = dot(modelview[1], vec4(vertexPos, 1.));
    viewSpacePos.z = dot(modelview[2], vec4(vertexPos, 1.));

    // use view space position in calculations
    ...

    // transform from view space to clip space
    gl_Position = viewSpacePos.xyz * projection;
    gl_Position.z += 1.0;
}
```

## 4.7. “Const” Data in Shaders

If used correctly the `const` keyword can provide a significant performance boost. For example, a shader that declares a `const` array outside of the `main()` block can perform significantly better than the same shader with the array not marked as `const`, even if the array could be treated as such. Another example would be the use of a `const` value to reference an array member. In this example, if the value is `const` the Graphics Core can know ahead of time that the number will not change and data can be pre-fetched prior to the shader being ran.

## 5. Optimizing Specific Techniques

### 5.1. Multiple Render Targets (Series6 only)

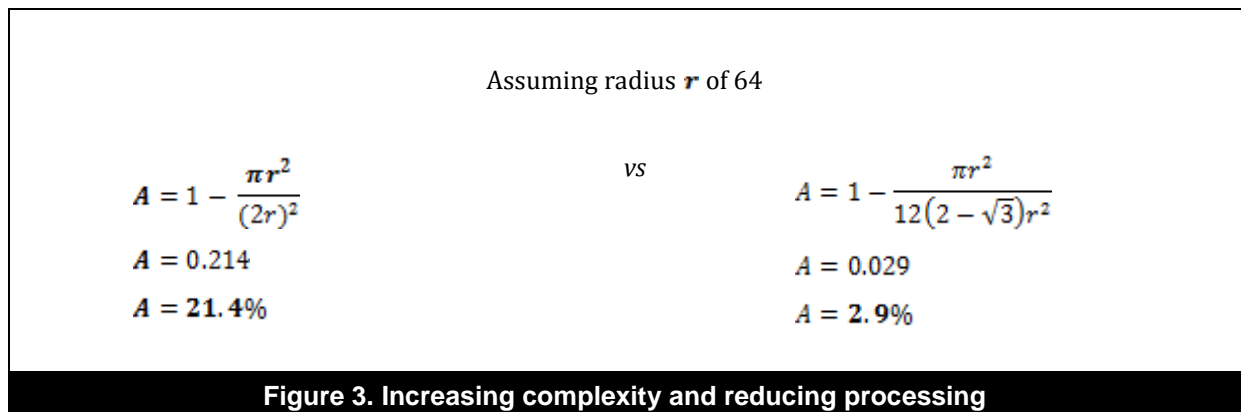
Multiple Render Targets (MRTs) are available in a variety of APIs, and are supported on PowerVR hardware from Series6 onwards. By using MRTs properly developers can take advantage of the tile-based architecture of PowerVR hardware, keeping all render targets entirely on-chip for a significant performance boost. In order to benefit from this feature the combined bit rate of all MRTs should be no more than 128bits per pixel.

### 5.2. Efficient Sprite Rendering

Rendering sprites efficiently may seem like a trivial exercise. However, without careful consideration an application may be unresponsive and sluggish due to poor graphics performance. Traditional sprite render tends to see textures drawn, using alpha blending, on to quads. These quads will consist of large areas of alpha, either full alpha, or partial alpha. Areas of full alpha are traditionally discarded using either the `discard` keyword or alpha testing, while areas of partial alpha undergo blending. Both of these have some form of impact on performance versus fully opaque objects meaning that a large number of sprites being drawn inefficiently can seriously harm performance.

The `discard` keyword (see “PowerVR Performance Recommendations: The Golden Rules”) should be avoided in favour of the much faster alpha blending. Even when favouring alpha blending, performance can still be affected if there are a large number of sprites. One method to minimise the impact of several layers of blended sprites is to increase the geometry complexity of the sprites in order to reduce the amount of wasted transparent fragments. For example, if a sprite is circular in shape and is rendered using the most optimal fitting quad, 22% of the fragments processed are redundant. Significant performance improvements can be gained by reducing the wasted transparency by increasing geometry complexity.

PowerVR hardware has excellent vertex processing capabilities and is designed to handle large amounts of geometry data, far in excess of what is present in most sprite based applications. As such, increasing complexity should have minimal performance impact and any impact this may have is most likely outweighed by the savings of rendering less transparency. If we increase the complexity of the previous case of a perfectly fitting quad around a circular sprite to that of a dodecagon (twelve sided polygon) we can reduce the amount of wasted fragment processing to just 3%.



## 6. Contact Details

For further support, visit our forum:

<http://forum.imgtec.com>

Or file a ticket in our support system:

<https://pvrsupport.imgtec.com>

To learn more about our PowerVR Graphics SDK and Insider programme, please visit:

<http://www.powervrinsider.com>

For general enquiries, please visit our website:

<http://imgtec.com/corporate/contactus.asp>

Imagination Technologies, the Imagination Technologies logo, AMA, Codescape, Enigma, IMGworks, I2P, PowerVR, PURE, PURE Digital, MeOS, Meta, MBX, MTX, PDP, SGX, UCC, USSE, VXD and VXE are trademarks or registered trademarks of Imagination Technologies Limited. All other logos, products, trademarks and registered trademarks are the property of their respective owners.