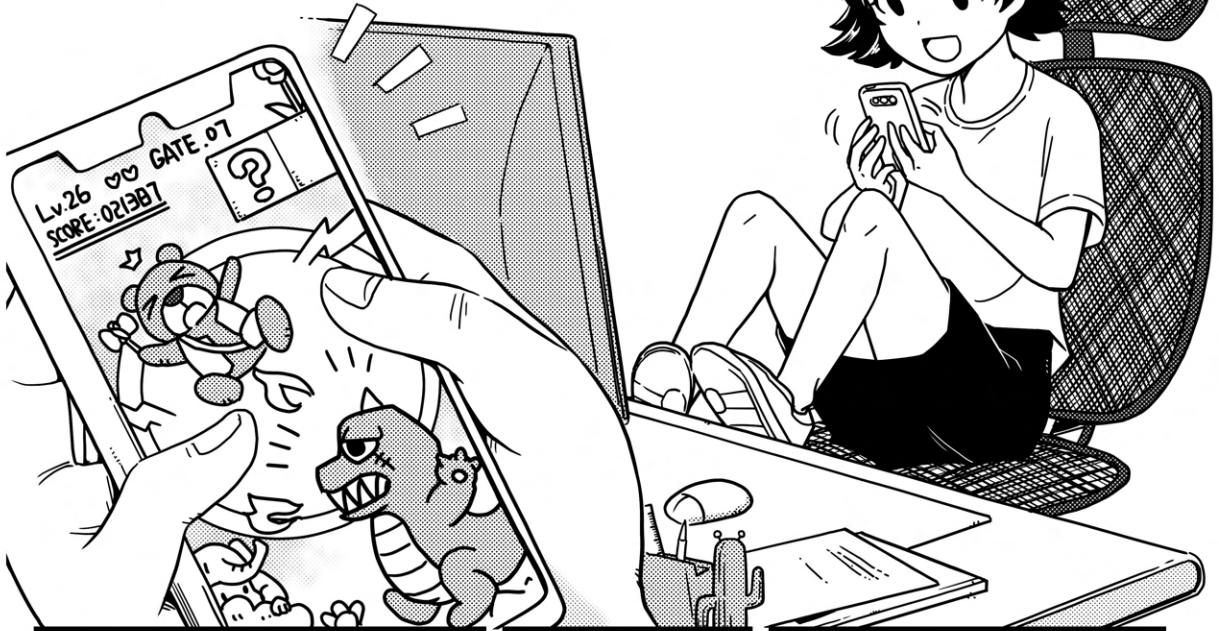


How does a Mobile GPU Work?

Comic: Ikaridon Yu

Playing games on mobile is so convenient.




The current mobile GPU is so powerful.

I can just port my PC game to mobile directly now.

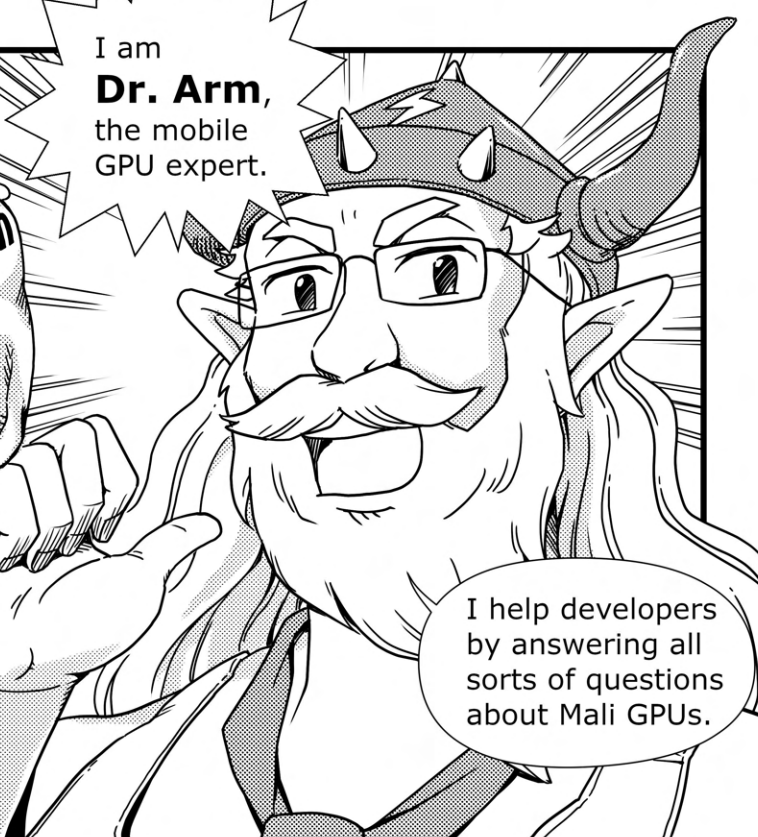
Hold on!!!

Although the mobile and PC GPUs do the same job, directly porting from the PC is not a good idea!

W.....Who are you?



Sorry, I forgot the intro.



I am **Dr. Arm**, the mobile GPU expert.


I help developers by answering all sorts of questions about Mali GPUs.




suspicious guy....

Why did you say that the direct porting is not a good idea?

Although the mobile GPU uses the same API as the PC GPU, the architectures of the two GPUs are quite different.



Hmm... but I still can't imagine the difference even though you said so...



Not a problem, let me take you inside a GPU.

LET'S GO!



WOW!

arm

developer.arm.com/graphics - 2 -



Here is the inside of a PC GPU.



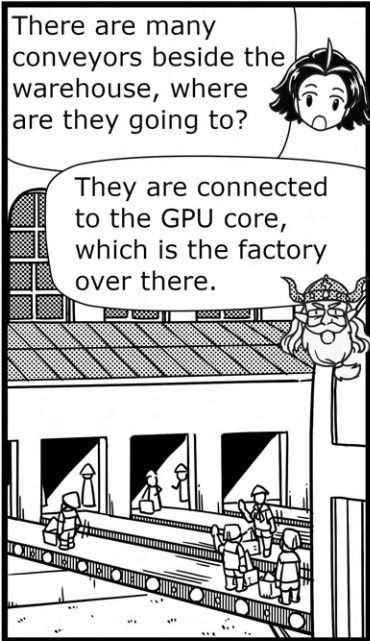
Wow! That is a huge power station.

That's right. That's why the PC GPU doesn't need to worry about power consumption.



Did you see a big warehouse over there?

That is **video memory**.
All data that needs to be rendered is stored inside it.

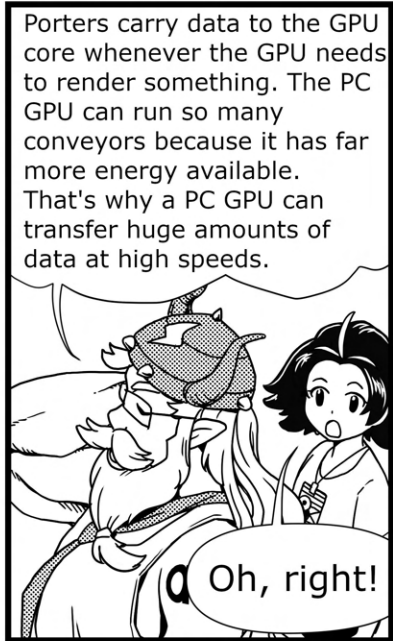


There are many conveyors beside the warehouse, where are they going to?

They are connected to the GPU core, which is the factory over there.

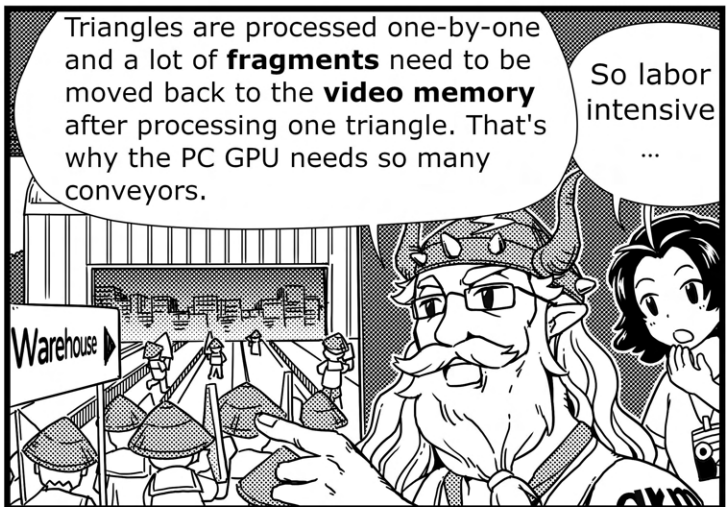
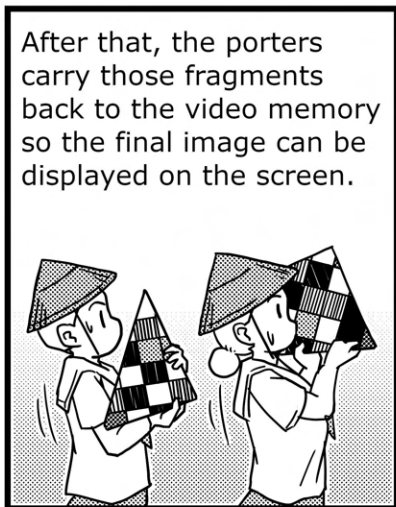
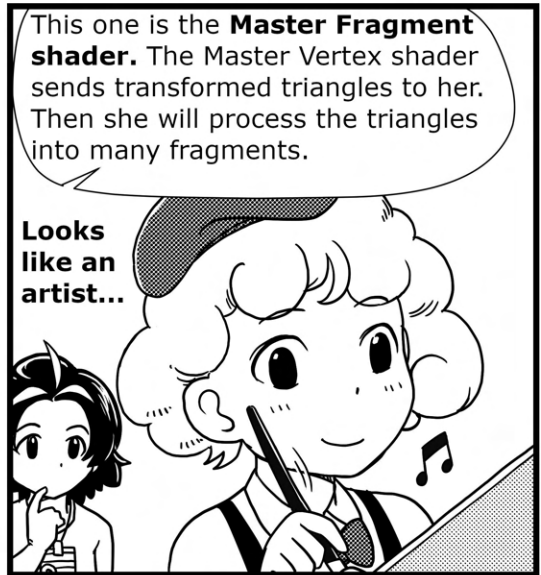
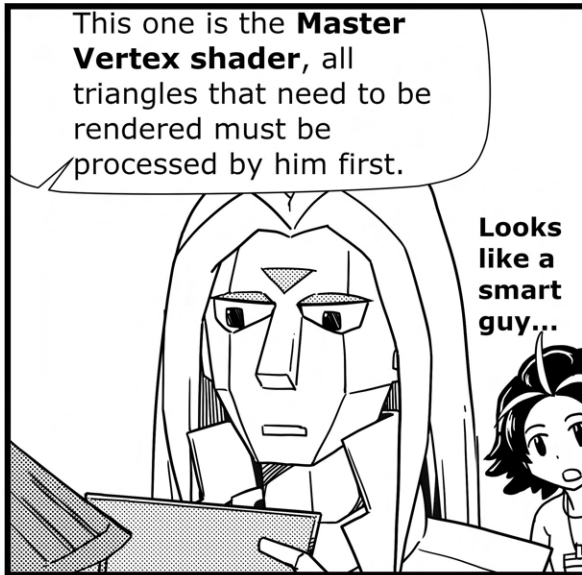
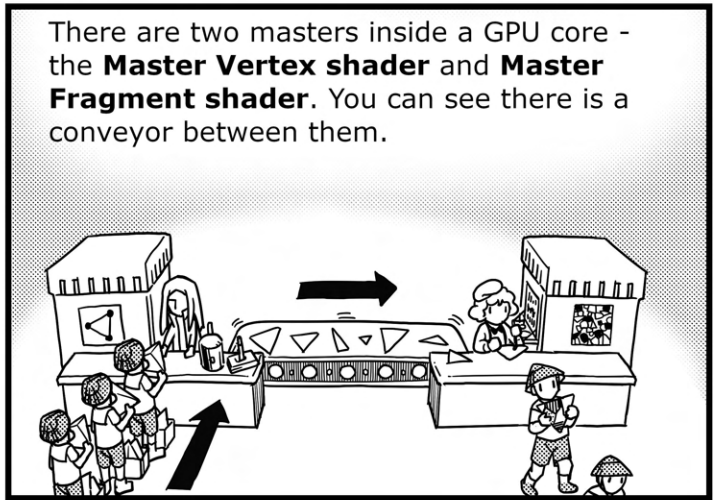
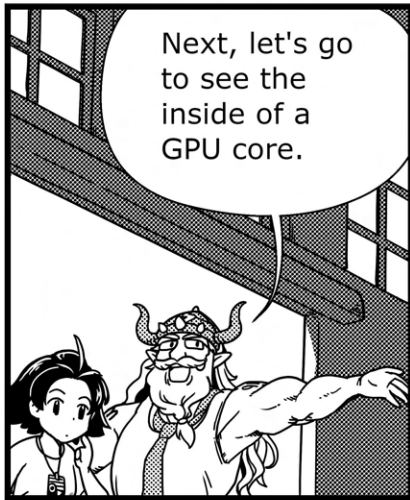


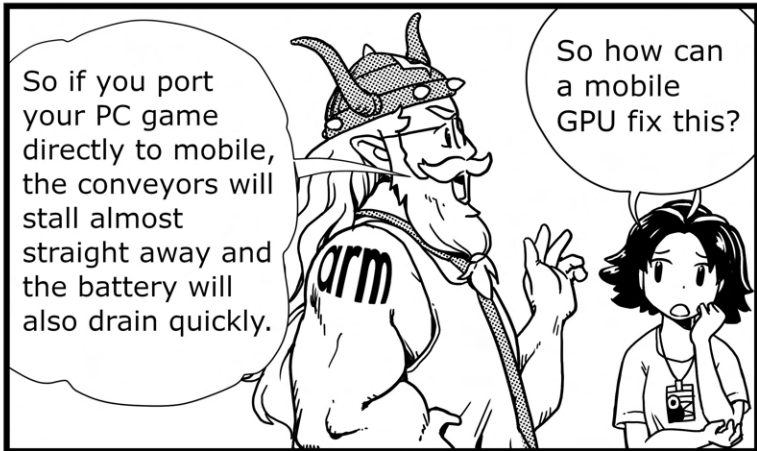
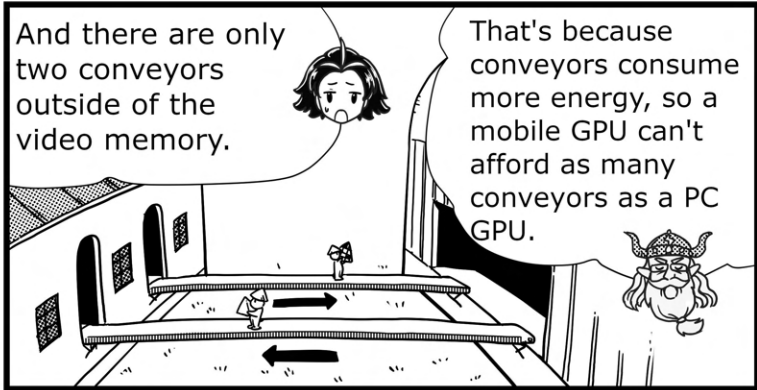
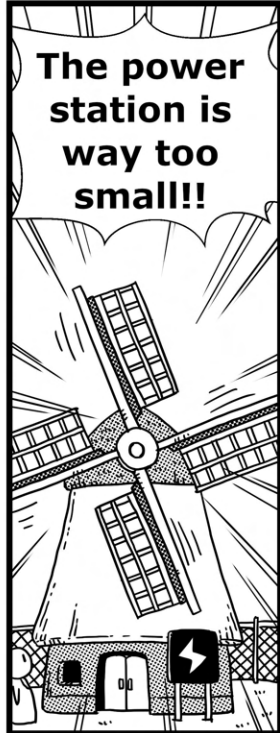
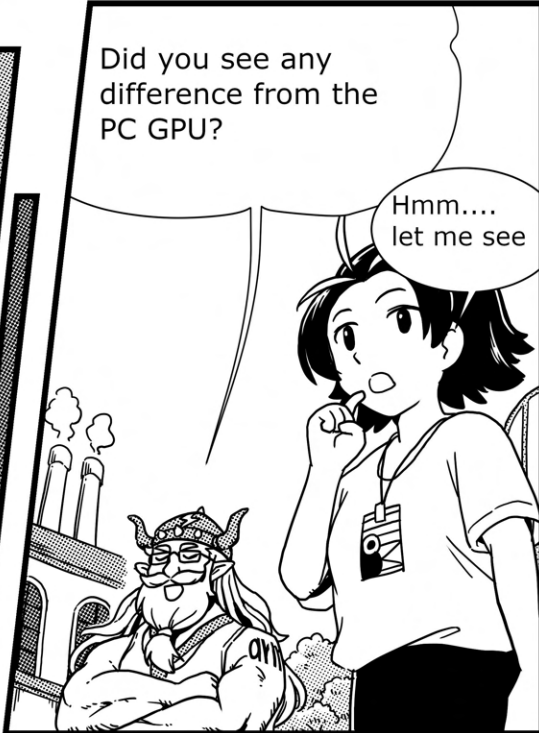
There are so many porters carrying data over the conveyors!



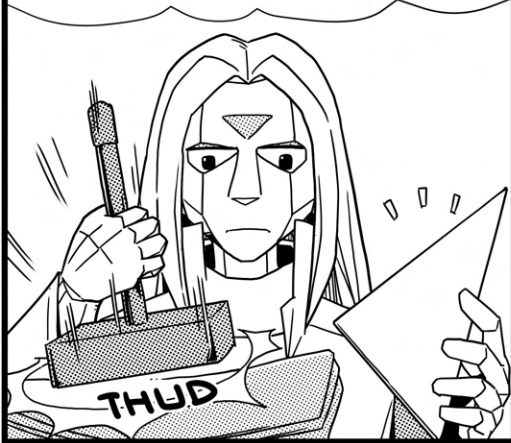
Porters carry data to the GPU core whenever the GPU needs to render something. The PC GPU can run so many conveyors because it has far more energy available. That's why a PC GPU can transfer huge amounts of data at high speeds.

Oh, right!



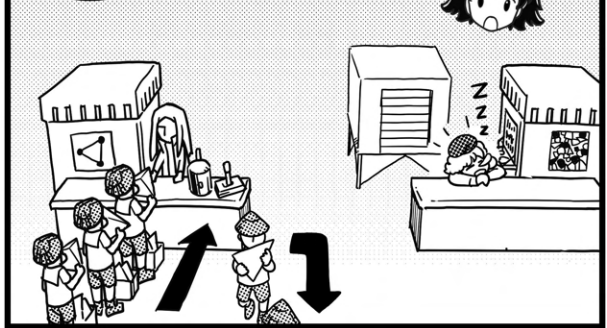


Let's take a close look at the mobile GPU core. It has the **vertex shader** to process triangles, but...



The transformed triangles will not pass directly to the **fragment shader**. Instead, the triangles will be transferred to the **video memory**.

What?! So what should the **fragment shader** do then?



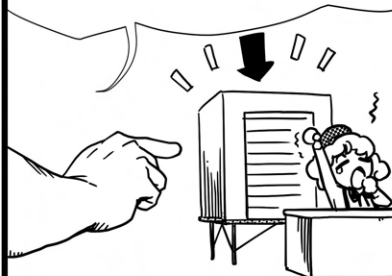
The **Fragment shader** will get the transformed triangles from the video memory then.



But would it increase the size of the data transfer?



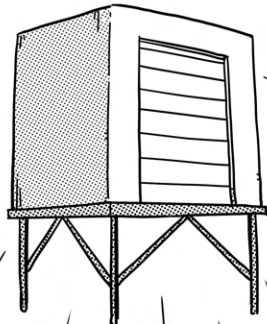
Exactly! Did you see a small warehouse besides the **fragment shader**?



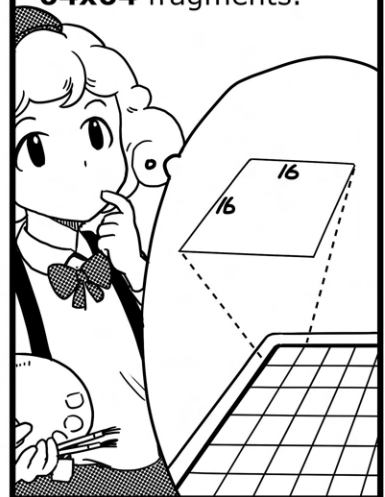
Yes, I didn't see such a thing at the PC GPU core.

It's **tile memory** a special design for mobile GPUs.

Tile



On a mobile GPU, the screen is split into many **tiles**. Each **tile** contains **16x16 ~ 64x64** fragments.



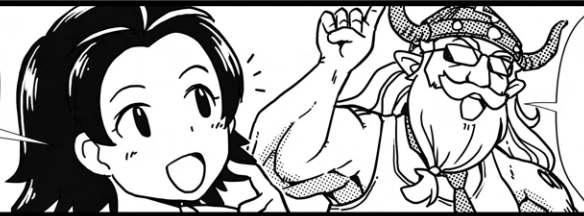
The **fragment shader** processes one **tile** at a time and then processes all triangles in the tile at once. So the frame buffer data just needs to be transferred to the **tile memory** at the beginning of the tile processing.



Then the **fragment shader** can directly access the data in the **tile memory** and process all the triangles at once.



I got it now. Although the data transfer for vertex is increased, the tile design saves more of the data transfer for fragment.



Correct. And there is another advantage to using the tile for rendering...

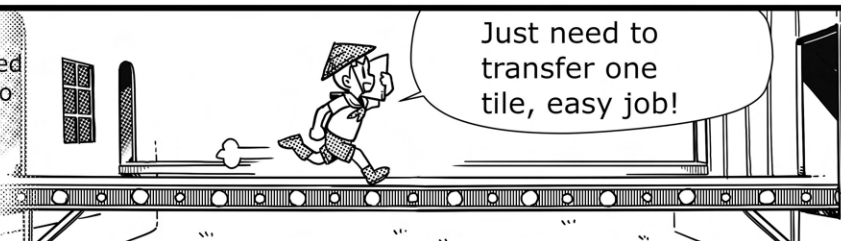
If the GPU finds any triangle that is **occluded** by other triangles, those occluded triangles will be discarded without rendering.



Cool, so using the tile can save both memory bandwidth and power consumption. What a smart design.



The whole tile will be transferred back to the video memory when all triangles in this tile are processed.



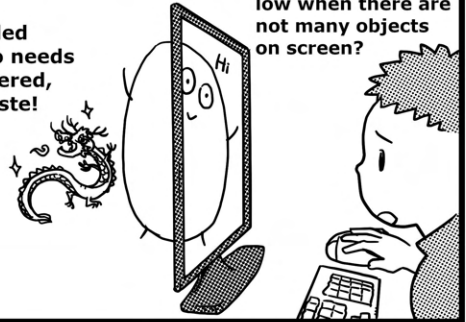
Because the number of **tiles** on a screen is fixed, the size of the data transfer is also fixed.

My size is fixed

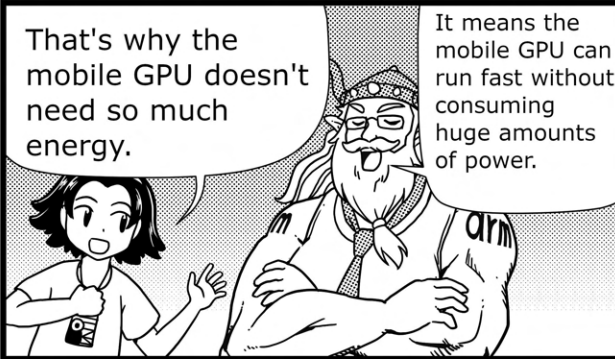


But on a PC GPU, the data transfer size is not fixed. It varies depending on the number of triangles that are rendered.

The occluded object also needs to be rendered, what a waste!

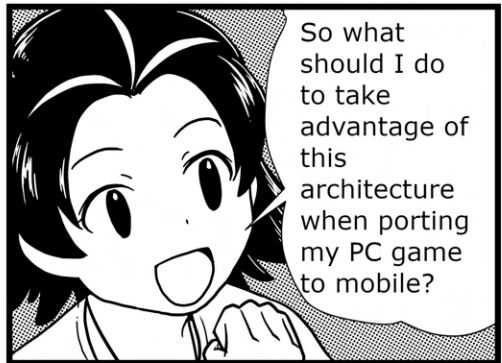


That's why the mobile GPU doesn't need so much energy.



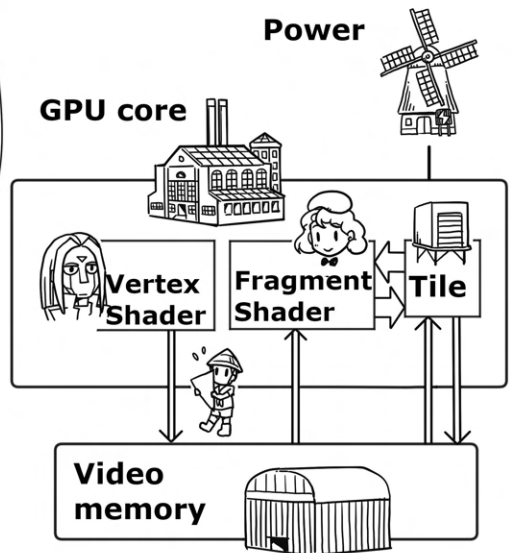
It means the mobile GPU can run fast without consuming huge amounts of power.

So what should I do to take advantage of this architecture when porting my PC game to mobile?



Good question. But we have learned enough today.

Let's review what we've learned so far and I will show you how to do that in the next episode.



Go to the Arm developer site for more detailed information :

<https://developer.arm.com/solutions/-graphics-and-gaming/developer-guides/-learn-the-basics/tile-based-rendering>