**Nanoelectronics: Devices and Materials**
**Prof. Navakanta Bhat**
**Centre for Nano Science and Engineering**
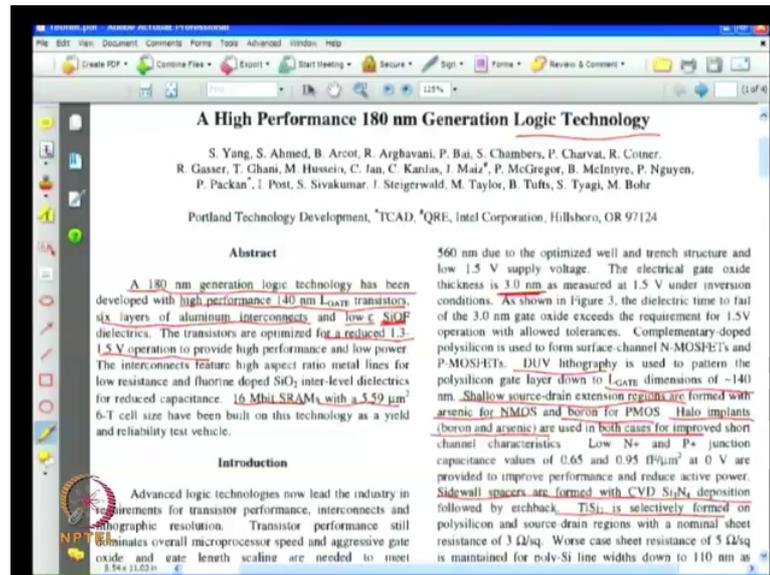**Indian Institute of Science, Bangalore**

**Lecture - 11**
**Industrial CMOS Technology**

Today, what we will do is; go through some representative CMOS technology. And you will recognize that all that we have studied so far is very relevant and it is in fact, being implemented in state of our CMOS technology.

So, what I thought I will do is; walk you through some representative notes. We will start with 180 nanometer generation CMOS technology and then progressively look at scaling; look at 130 nanometer, 90, 65, 45, 32; each one being 0.7 X of the previous technology. So, let us then start with this; what I will do here is essentially I will make use of these papers, which are published in conference called IDM; which is a annual conference called International electron Devices Meeting; it is a premier conference in electron devices. And any new device technology, which gets developed, gets reported in these conferences for the first time.

So, what I have done is that; I have taken a series of papers published in past several years. We will go through these papers one by one and of course, we will not really read the entire paper, but I will highlight a few important things to draw your attention to some of the very important technological implementations.

Let us look at this paper here now; which is essentially A High Performance 180 nanometer Generation Logic Technology. And just to be consistent; 180 nanometer technology and all subsequent technologies of course, will be implemented by various companies, various foundries. Just I have picked papers published from Intel; just without any bias to any company, but just to tell you a representative technology and also with respect to Dell technology, Intel technology; we will see how the technology has progressed.

So, this is essentially a 180 nanometer Generation Logic Technology, if you see here 180 nanometer Generation Logic Technology; has been developed with high performance 140 nanometer gate length transistor. You see, I had told you that there is so called technology node, but the gate lengths of the transistor could be less than the technology node. So, this technology node is 180 nanometer. However, the gate length of the transistor is actually 140 nanometer. It has 6 layers of aluminum interconnects, the interconnect in this technology is still aluminum. However, later on you will see that they have migrated from aluminum to copper; in terms of the interconnect.

And the dielectric between interconnects, which is metal interconnect; if you recall, we have discussed this in our CMOS process flow; between one metal level to another metal level, you have an insulator. And that insulator is silicon oxide; except that it is little bit variant of silicon oxide; it is fluorine doped silicon oxide. The idea of doping it with

fluorine is to reduce the dielectric constant of silicon oxide; that is why they say it is a low Epsilon; Epsilon is a dielectric constant.

Because in interconnects you want to minimize the parasitic capacitances, parasitic resistances. And hence you need to minimize the capacitances, this is opposite of high k gate dielectric, these are low k gate dielectrics. So, you minimize the interconnect delay and interconnect capacitances. The transistors are optimized for a reduced 1.3 to 1.5 volt operation you see the transistors operate at low voltage; at 1.5 volt maximum, not beyond that.

The interconnects which are high aspect ratio, metal lines and all that and using this technology they have demonstrated a 16 bit SRAM; with certain memory cell size and this is sort a demonstration of the technology. What they are showing here is that; not only they have developed the transistors, but they have put together transistor to make a circuit. And in this case, the circuit is static random access memory; which is a representative circuit.

Eventually in this technology is a logic technology; you see here logic technology, not a memory technology. Logic technology will be used to build microprocessors, which are very high and logic chips, but microprocessors on chips will have cache memory. So, these SRAMs will be used for cache memories in the microprocessor, these are not going to be standalone SRAM chips to be sold.

So, we have been able to demonstrate that SRAM circuit can also be implemented using this technology.

(Refer Slide Time: 05:10)



So, let us then just go through some key highlights here and they give the introduction; why do you need scaling and all that; that is very typical.

(Refer Slide Time: 05:23)



Now, let us come to the transistor; so, the transistor is sort of illustrated in this figure; you have CMOS technology, you have N channel and P channel transistor. Figure one illustrates the structure of the MOS transistor and isolation used in this technology. The isolation as we know is a shallow trench isolation; correct now tello coat isolation. You have made a trench inside the silicon and that isolates the neighboring transistors. Start

with a P minus P plus epitaxial silicon wafers, followed by the formation of shallow trench isolation, N wells are formed with deep phosphorous and shallow arsenic implants; they have not used in this particular case antimony and indium yet.

But none the less they use phosphorous and arsenic combination for N well; while P wells are formed with boron implants that are how they are forming these wells. And then they also have the electrical gate oxide thickness is; 3 nanometer. Gate oxide what they show here is a electrical thickness, when we talk of C V characterization subsequently. From next class, we will see what is the difference between physical thickness and electrical thickness?

Electrical thickness is always little more than physical thickness. So, if they are saying 3 nanometer is their electrical thickness, their physical thickness could be something like 2 nanometer; less than 3 nanometer for sure. And they have ensured that; this dielectric has sufficient reliability as they will demonstrate and their gate lengths have been patterned using deep UV lithography; it is an optical lithography essentially.

Shallow source drain extensions; remember all that we talked about in terms of source drain engineering, they have used shallow source drain extensions with arsenic for NMOS, boron for PMOS; these are source drain extension for N channel and P channel transistors. Halo implants, boron and arsenic are used in both cases; for improved short channel characteristics, low N plus P plus junction capacitance values; that is fine.

Sidewall spacers are formed using CVD nitride; remember the nitrides spacer that we talked about, that is in fact; used to isolate deep sourced extension from your shallow extensions, they make use of that. They use silicide; titanium silicide, remember the parasitic region that we talked about; after the gating of the transistor, there is a contact and in between there is a silicon and if you convert it into a silicide, you get lower resistivity and in this technology they are using titanium silicide.
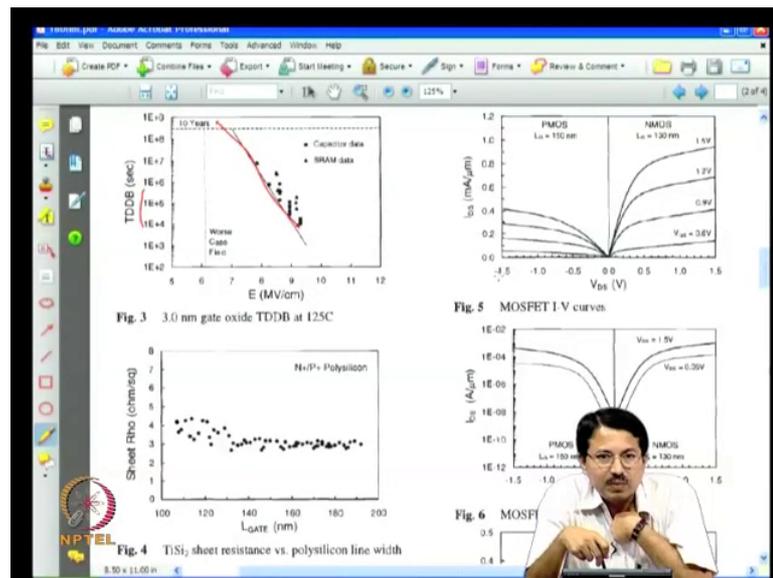
So, that the typical on; and source drain regions with a normal sheet resistance of 3 ohms per square; that is the kind of sheet resistance that you can get, which is a very low sheet resistance by the way. So, hence this is how the cross section would look you see, there is a N channel transistor in a P well, there is a deep source drain, there is a shallow extension, there is a spacer; which separates the two and there is a silicide here. And

similarly silicide will be on poly, as well as on this junction and this is a PMOS transistor; complimentary metal oxide silicon technology.

So, this is how you would define the transistor and then you will have to characterize a lot of things. You will have to make sure that the isolation is sufficient, for that you need to characterize the breakdown voltage of these; that is between this and this. If you apply 2 voltage; it should not break down, the breakdown voltage should be more like 5, 10 volt. So, for all practical purposes these two are perfectly isolated, so there is essentially demonstrating all these characteristics.

You can go back and go through these papers, these papers are certainly available if you are an I EEE member or if your institute is subscribing I EEE journals; you will certainly have this listed in the I EEE publication.

(Refer Slide Time: 09:31)



Another important aspect that we have discussed; all these reliability issues: remember I said when you do the gate oxide reliability; there is something called TDDB; which is time dependent dielectric breakdown. Even though you are going to operate your transistor at very low electric field, such as may be in this case 6 mega volt per centimeter; what you do is that, you apply large electric field and characterize the breakdown times and based on that you sort of extrapolate it, and say that in the real use condition certainly this will be reliable up to 10 years; they have done a similar thing, what we have discussed; these are typical reliability characterizations.

And here they show typical N channel and P channel characteristics, this is drain current versus drain voltage; the so called output characteristics, which are well behaved; going up to 1.5 volt for NMOS and minus 1.5 volt for P channel transistor. Of course, here N current is more than the P current because of the electron mobility being higher than the whole mobility.

(Refer Slide Time: 10:54)



This is sub threshold characteristics your threshold voltage is about 0.5 volt or so, and below that this is IDS and log scale and VGS in linear scale; it is a semilog plot really and this really is a sub threshold slope, as we have already discussed. And again notice that; depending on whether you have 0.05 volt here or 1.5 volt here, you will have different characteristics; that are essentially due to double effect drain induced by real lowering.

If you have larger drain voltage for the same gate voltage, you get larger current because your drain electric field will influence the injection at the source side; especially in these nanometric dimension transistors. So, this is just to show that; your polysilicon gate because of the silicidation has reasonably low resistivity.

Sheet resistance is the order of 3 to 4 ohm per square, which is fair good. These are typical Vt versus length plots, remember what our discussion; there is a reverse short channel effect. I told you whenever you use a halo; invariably you will have reverse short channel effects. And then usual Vt roll off that you will have in this region here this; Vt

is coming down. This is Vt measured at 1.5 volt on the drain and this is Vt measured at 0.05 volt on the drain and this difference is essentially DIBL.

Vt at low voltage on the drain is always higher than Vt at high voltage on the drain; your DIBL can be characterized based on that.

(Refer Slide Time: 12:32)



So, they have NMOS and PMOS devices which have a sub threshold slope of 90 millivolt per decade. Remember 60 is the best you can get, but you will never get 60; because you have 1 plus CD over C ox term, which is always more than 1. So, because of that it will be more than 60; so, they have been able to achieve 90 millivolt per decade. And their off currents are of the order of 3 nano ampere per micrometer, we always specify off current in terms of ampere per unit width; depending on how many microns your transistor width is you multiply that, you get the off state current.

And your on state currents; saturation drive currents are 0.94 milliampere for NMOS and 0.42 milliampere for PMOS. So, in other words your on current to off current ratio is almost 10 power 6, if it were to be 10 power 6; then it would have been 3 milliampere here, but it is not 3 milliampere; it is close to 1 milliampere, which is reasonably good on current to off current ratio. This is NMOS, this is PMOS same story; reverse short channel effect and typical Vt roll off; the Vts are negative.

Again lower drain voltage, higher drain voltage; short channel, threshold, voltage, roll off are shown in these figures; threshold voltages are 1.5 volt; drain bias are 0.3 volt for NMOS and minus 0.24 volt for PMOS. And these results are of course, better than any previously published bulk or SOI devices; they are trying to say that, they have been able to come up with much better devices.

Subsequently, they have also made use of these transistors and dealt what is called ring oscillator and from that they extract inverter gate delay. If you were to build a simple inverter, a not gate; what is the gate propagation delay? You have a input transition from low to high, how long does it take for output to go from high to low? That is your propagation delay and they have been able to get the propagation delay and typically, it is very difficult to measure propagation delays of the order of these propagation delays that they have are of the order of few picoseconds.

You do not construct just one single inverter; because if you want to construct one single inverter, you will have to measure the time resolution of the order of picoseconds. What you do is that; you cascade the inverters, a large number of inverters are in series; very interestingly, if you have a odd number of inverters in series and connect the output of the last inverter back to the input of the first inverter.

Create a feedback circuit; that starts working as an oscillator that is why it is called a ring oscillator. Ring meaning, you have created a inverter ring; you need to have an odd number of inverter; if you have even number of inverter, it will go to a steady state; input will also be a study, output will also be steady, but if you have a odd number of inverter; you can do this exercise yourself; output will try to change the input and it will keep changing; it will become a oscillator circuit.

So, you do this oscillator; you measure the oscillators time period and you know how many stages of inverters you have and you back calculate; the delay per inverter and that is what is done typically. And they do it for unloaded ring oscillator meaning there is no extra load on the oscillator. The only load is the subsequent inverter stage; operating at 1.3 and 1.5 at room temperature; they have been able to get very low delays of the ring oscillators.

(Refer Slide Time: 16:39)



We will come to this in a minute, but this is typically what is plotted; gate delay as a function of gate length. If you have a long transistor, your delays are also more; that is why we are scaling; I mean the reason why we scale is that, the circuit start operating faster as we have discussed in the very beginning. As you start shrinking this gate length, the gate length is decreasing here, the delay is also decreasing.

And the nominal gate length for this technology is of the order of 150 nanometer and that 150 nanometer, your gate delays are of the order of ten picoseconds. Again, if you operate this at lower voltage; your delay is little longer, little more and if you apply higher voltage for this circuit; not very high, I mean if you apply very high voltage; it will break down; then you can decrease the delay a little further.

One good way of bench marking the transistors is to really plot this; so, called DC performance metric; what is DC performance metric? This is a XY plot with on current on the X axis and off current on the Y axis. And the requirement for a good transistor is that; for any given off current, whatever that off current is that; let us say 10 nanometer off current.

I want to be furthest right of this curve, if you have two technologies for example, if I have a new technology; which looks like this. This is much better transistor because for the same leakage current, it is giving you much much better on state current; on state current is much larger. Here on to off current ratio is huge, so that it will be less leaky; it

will not consume large static power, at the same time it will be very fast. So, it is essentially a 2 D plot of off current versus on current and you want to really create your transistor to really go along this; towards the right.

And of course, NMOS tends to be for the same leakage current; to the right of PMOS, simply because you have much better electron hole mobility; compared to hole mobility; but what they are trying to do here is that, the filled simple they say this work and their some previous works; from references, these are empty symbols. So, what they have plotted here is that; these empty symbols are to the left of; all these are PMOS transistors from previously published papers.

And all these empty symbols here are NMOS transistors from previously published papers. And all they are trying to say is that; compared to previously published PMOS devices; these are much better. Because for the same off current, your curve is to the right; you get much higher on current and that is what you want to generate. The way you generate this curve by the way; if you want to generate such a curve, you cannot generate this curve with a one gate length transistor.

What you do is that; using the same technology, you print transistors of different gate lengths; gate lengths from 130 nanometer, 150 nanometer, 200 nanometer, 300 nanometer so on and so forth. As your gate length is increasing, your leakage current decreases, your on current also decreases. So, all these points here; you see correspond to shorter gate length; because these are large on current and large off current. All these points here correspond to longer gate length pitch channel transistor and that is how you generate series of these points by characterizing transistors of different dimensions.

And you plot a universal curve which is a off current versus on current curve. Similarly, for N channel transistors; these are the transistors which have very high on current and off current, which are coming from very short channel 130 nanometer kind of transistor. Whereas, along this these are longer channel transistors and that is how this is generated and they do interconnects. Interconnect is essentially aluminum in this case they have not yet used copper.

But the only thing that they have done different in interconnect is to use fluorinated silicon oxide as I mentioned earlier here.

You see fluorine is added to a silicon oxide; to reduce dielectric constant and improve interconnect performance. The use of SiOF as an inter level dielectric, reduces the dielectric constant to 3.55 compared to 4.1 for undoped silicon oxide; remember silicon oxide, the dielectric constant is around 44.1 and whereas, that comes down to 3.55.

You see thermal oxide, if oxide is grown at let us say high temperature; 1000 degree centigrade then its dielectric constant is more like 3.9; that is what we used earlier when we are talking of gate dielectric; 3.9, we round it off to 4. But here we say, 4.1 because this oxide is not thermally grown, it is a deposited oxide. Because you already have a metal line and top of metal line, you have to put oxide; so there is no silicon to grow oxide, you will have to do a chemical vapor deposition of the oxide.

And invariably these oxides which are done at low temperature, because you already have a metal; tend to have little higher dielectric constant and hence the dielectric constant is of the order of 4.1; after fluorine doping, it comes all the way down to 3.55.

(Refer Slide Time: 23:00)



So, they actually build a SRAM cell; this is the top view of the SRAM cell where in you have 6 transistors, this is 6 transistors SRAM cell. This is their interconnect stack; starting from the transistor at the very bottom, these are the contacts; metal 1 via 1; metal 2 via 2; metal 3, via 4, metal 4 and so on and so forth; going up to 6 levels of metals; in this case.

(Refer Slide Time: 23:34)



And using all these, they have been able to show a working SRAM and this is what they are sort of summarizing; 180 nanometer generation logic technology has been developed

and demonstrated with high performance, reduce power transistors; aluminum interconnects with low Epsilon silicon; oxide with fluorine doping dielectrics are used to meet interconnect density and performance requirement.

The technology yield and performance capabilities have been demonstrated on a 16 megabit SRAM, which operates at greater than 900 mega hertz frequency. So, eventually if you have to build a microprocessor; if the cache has to respond at that frequency, it would respond at that frequency.

So, they have been able to start from the basic semiconductor processes; using all the process integration, now build the transistors, demonstrate a simple circuit like ring oscillator, characterize the gate delay and take it forward to make an SRAM cell and show that the SRAM is working. So, this is really a demonstration of a technology, but still this technology is still not; after this it may take another year for you to make a real product. Because your real product is going to be a microprocessor subsequently, then you will have to make a microprocessor, make sure that yield is good.

(Refer Slide Time: 24:59)



And that requires lot of optimization; so, let us now look at the subsequent technology which is the 130 nanometer generation logic. Again logic technology because Intel does not make standalone memories, the only product they offer and that is; their essentially lifeline is microprocessors, very high end microprocessors. And hence this is the logic technology, not a memory technology again.

Now 70 nanometer transistors; although it is a 130 nanometer technology; you see the gate length is 70 nanometer. They have dual Vt transistor; remember something that we discussed, there is a leakage power problem, you need to really make sure that you minimize the leakage power; at the same time you need high performance. So, what you do? It turns out there are only is few critical parts, in any complicated chip.

It is important to make the transistors in the critical part very fast transistors and hence make only those transistors a low Vt transistors and others could be high Vt transistors. As a result of that you get very fast chip as well as low power consuming chip. So, in other words here; you have two flavors of N channel transistors and two flavors of P channel transistors. One is called the low Vt transistor and the other one is called a high Vt transistor.

And hence it is a dual Vt transistor technology and what else? 6 layers copper interconnects, a big departure from aluminum to copper. Because this is the only way to scale the interconnects, very briefly may be we will look at the issues and interconnect scaling; now let us get started what they have here. A leading edge 130 nanometer generation logic technology with 6 layers of dual damascene copper interconnect. Dual damascene is a particular process sequence which is used to realize copper interconnects.

We will not really go over the details of that process sequence; dual Vt transistors are employed with 1.5 nanometer thick gate oxide, the gate oxide thickness as gone down compared to what it was in the previous generation; operating at 1.3 volt; previously it was 1.5 maximum; now it is 1.3; you see. High Vt transistors have drive currents of 1 milliampere and 0.5 milliampere per micron; for NMOS and PMOS.

While low Vt transistors have 1.17 for NMOS and 0.6 for PMOS respectively; so two flavors of transistors. Again they have been able to make an SRAM cell, demonstrate a SRAM cell; demonstrate a technology on a 18 bit SRAM.

The typical process flow is the same thing; except that you have a two dual Vt; meaning you will have to create two different N wells, two different P wells. So, that doping concentration is little different in two N wells and little different in two P wells and that is how you will get dual Vt technology; for N and P, which means more photolithography step, little more complex technology. But that is worth it; only if you do that, you would be able to get the best of it.

So various technological features; this is how a transistor looks; N channel transmission electron, micrograph of showing 70 nanometer gate length; source drain, ultra thin gate electrode, polysilicons, silicide at the top and so on and so forth.

Or else will remain identical; you will use halo; you will use shallow extension except that you may change the dose; because the transistor gate length is going down especially halo dose, you may want to manipulate. Shallow source drain extension regions are formed with arsenic for NMOS, boron for PMOS. Boron and arsenic halos are used; this will remain same, silicon nitride is used for the spacer; correct? Followed by etch-back; the silicide here is cobalt silicide as a (Refer Time: 29:28) titanium silicide.

Cobalt silicide gives you much better, lower resistance compared to titanium silicide. Again the demonstrator well behaved NMOS and PMOS; it is the story looks similar right, you will see similar graphs here except that now you have a high Vt and low Vt; two flavors.

And sub threshold plot, drain current in a log scale versus gate voltage; for low drain voltage and high drain voltage. And the DIBL for the 70 nanometer NMOS device is measured to be 100 millivolt per volt that is how we specify the DIBL. What is the difference in Vt per unit voltage; on the drain; unit voltage change on the drain? And that is how they are characterizing the DIBL.

(Refer Slide Time: 30:26)



And this is again the typical off current versus on current and they are sort of saying that; this again two flavors unfilled and filled, high Vt and low Vt; there two flavors of PMOS and two flavors of N channel transistors; these are some important matrix.

(Refer Slide Time: 30:46)



Such as VDD, gate length, oxide thickness; they are comparing 180 nanometer generation, which is a previous paper that we discussed and this generation; supply voltage from 1.5 to 1.3, gate length from 130 to 70, oxide thickness from 2 nanometer to 1.5 nanometer and so and so forth.

(Refer Slide Time: 31:13)



You see off current has increased a little bit; from 3 to 10 nano ampere per micron. Typical threshold voltage role of; threshold voltage versus gate length, N channel P

channel; N channel positive and P channel negativity; low drain voltage, high drain voltage; this is because of DIBL. This is now high Vt and this is the low Vt transistor.

So, again they characterize delay; they build ring oscillator. Now they have depicted in slightly different fashion, in the previous paper; you saw delay plotted against gate length. Here, delay is plotted against off state current of a transistor; remember off state current is inversely related to gate length. In other words, if you have to plot gate length here; this would have been a lower gate length, may be 100 nanometer here, may be 200 nanometer here, may be 500 nanometer here.

500 nanometer will also have; I am sorry, I got it the other way around; this would be 100 nanometer, let me see if I can delete this; let me just take it away; that is ok.

(Refer Slide Time: 32:40)



So, this needs to be sort of reverse; 100 nanometer is out here, 100 nanometer also results in large leakage current; longer channel length also results in small leakage current correct, I mean let me write it here so, that becomes very clear.

(Refer Slide Time: 33:09)



So, this is 100 nanometer transistor, let us say 100 or 100; 70 nanometer in this case, they have gone down to 70 nanometer. So, this would be something like 70 nanometer out here and this may be 100 nanometer and this may be 200 nanometer. So, the gate lengths is a lowest here and the lowest gate length gives the lowest delay.

And when the gate length is increasing, the leakage current also decreases as is evident here; lower leakage current. Because you have larger gate length, the on current will also be lower; just as off current is decreasing, on current is also decreasing and hence you will get larger delays. It is essentially looking at this delay versus gate length or delay versus leakage current; one and the same essentially.

So, chip performance is increasingly limited by the RC delay of interconnect; now why is that?

Let us look at what is happening with scaling; what have we done over scaling? We have started from let us say technology of the order of the 250 nanometer, quarter micron technology 180 nanometer, 130; 90, 65, 45 and so on and so forth. And let us look at what happens to delay; we know that as I scale transistor, transistor become very efficient. In other words, if we are talking of gate delay or inverter delay; over the technology generation, inverter delay or gate delay has been decreasing and that is why in the first place we want to scale.

You build any gate, it could be a not gate and gate or what have you right. So, called gate delay decreases; gate delay decreases because remember? Gate delay depends on transistors and transistors with scaling become efficient, smaller transistors are always better and the so, called CV or I metric become better; as I starts scaling down the transistor. What is the story with interconnects? Interconnects, it turns out the story is completely opposite unfortunately.

(Refer Slide Time: 35:37)



What is an interconnect? Interconnect is a metal line, let us say this is metal 2 and there may be another neighboring metal 2 line here; this is another line running at the same level. And let us say; there is another metal one underneath M 1; between these two, you have inter level dielectric, two metal levels between that you have a dielectric; that is how you have insulated; this is metal 1.

What is scaling of interconnects? Scaling of interconnects is scaling the dimensions of the interconnects; just as scaling of the transistor, dimension has to come down and transistor distance has to come down. And same here with the interconnects, the interconnect; their dimension has to come down, what is the dimension? You see; this is the so called line width; we call it line width.

We print the metal line; metal line is a long metal line; how long it is; depends on where is this metal going in a chip from one location to the other location in a chip. This is the thickness; let us say t and this is length, let us normalize the length. Unfortunately, the chips are not becoming small; the chips are really becoming bigger and bigger. You would have thought that with scaling because you are trying to bring down the number of transistor, the chip should become smaller. That is why we say at 0.7 x; so, that the chip area comes down, transistor area comes down and all that. But what has happened over the years is that, with scaling we also put in more transistors in a new technology; we are not just happy with the performers gained that come just because of the mere scaling.

We want to increase a number of transistors very significantly compared to what it was earlier. So, if previously you had million transistor, you want to put 4 million transistor in a new chip, in a new technology. And hence you continue to use the same area or even bigger area. So, in other words if you are looking at a chip; the chip size from one technology to the other technology really does not scale; if anything it is becoming bigger.

The chips are becoming bigger; what is the interconnect? Interconnect is essentially a metal line which takes the signal from one location to the other location in a chip you see. The message here is that, the lengthwise it is not really helping us; it is really becoming as big. So, let us look at normalized length what happens with interconnects scaling; what is interconnect delay? Interconnect delay depends on R and C of the interconnect; this is the R C delay; due to interconnects.

Scaling, I need to scale the width, I need to scale t; if I go from this technology to a new technology of metal interconnect, in that the thickness has come down, the width has come down. What is the implication? Resistance goes up; because your cross section area for conduction is going down. So, per unit length; you see, if you were to look at ask the question what is resistance per unit length? What is rho L by A? Let us say I am interested in R per L; it is rho by A; how is A scaling?

A has t and W; as a result your interconnect resistance per unit length goes up as k square; which is not a good thing, this is not the end of the story. So, the message is that interconnect lines are becoming more resistant with scaling; not the end of the story, there is even more serious problem and what is that serious problem? You would have thought that at least because I am making interconnect smaller, the capacitance should go down.

Because capacitance should be due to; let us say, this metal line and another metal line; there is a coupling capacitance and the area for that coupling is t; I mean this that is determine by this width, that is coming down and hence it should go down.

Yes this so called parallel plate capacitance or inter metal capacitance indeed goes down, but there is a more serious issue and what is that? I am also decreasing the distance between the two metal lines in the new generation of the technology. So, the new

generation of the technology has M 2 lines; which were far apart, have been brought very close to each other.

So, there is a new capacitive component; which is intralayer capacitance. Earlier this distance was very large, you could have ignored that capacitance. Now this is increasing very dramatically, so net result is that with scaling your R is increasing, your C is also increasing; as per as interconnects are concerned.

C is increasing very dramatically because these interconnects lines are coming close to each other. And hence the R C delay of interconnect is increasing. So, if you were to look at interconnects, in older generation technology; 250 nanometer and beyond, interconnects were afterthought; you never bothered about interconnects. Their ideal metal lines takes signal from one location to other location, they do not contribute any further delay.

So, it was nonexistent compared to gate delay; your interconnect delay was very very low. But what has happened in the recent past is something like this; your gate delays are increasing, your interconnect delays are increasing. And hence interconnects are becoming bottle X; it is like just imagine your road network in the city, you have the fastest driving car. And that is transistor for you; that is a gate for you, but you have a very bad road network.

Then you are limited; not by the speed of the car, but you are limited how fast the roads will allow you to travel. And these are the roads; interconnects are the road that allows you to take signals from one location to other location; and that has become a serious problem. And then you see that in this technology; there are lot things that are being done; to minimize R and to minimize C.

In previous paper; we saw fluorine was doped to minimize capacitance and today I mean in this technology, we see aluminum is replaced with copper because copper has lower resistance compared to aluminum. RC delay becomes much better and that is exactly what is done here. And they have optimized this fluorinated oxide and that gives a k of about 3.6 and this is again the hierarchy of interconnects.

(Refer Slide Time: 43:06)



And they have used this to build; what they are showing here is that; if they use aluminum, you see with decreasing pitch, decreasing pitch is two metals are coming closer to each other. Your delays are increasing, if you are use aluminum; it would have gone up like this. Because you use copper, you could decrease it; at the same pitch, copper technology gives much lower delay compared to aluminum technology.

(Refer Slide Time: 43:42)



And that is what they have done and they have been able to make an SRAM and demonstrate that 18 bit SRAM is operating now at what; 1.6 mega hertz; much higher

frequency. So, actually it is not very clear here, I take back what I said; I mean the previous one was a lower density SRAM was at 900 mega hertz, but this is a higher density SRAM; that is operating at little lower frequency. As you start increasing the density, there is also an impact on the speed of the SRAM that you can build.

(Refer Slide Time: 44:38)



Whereas the ring oscillator; obviously, are much more efficient compared to the basic technology is much better. Previous generation was 10 picosecond; 10 to 13 picosecond, now it is a 7 pico second. So, you can clearly see a better technology and if you go down to 90 nanometer; now what has happened now? 7 layers, 6 layers have become 7 layers of copper; 90 nanometer has a 50 nanometer gate length transistors.

From 70, it came down to 50 nanometer now, it is using a strained silicon channel; we have not discussed about strained silicon channel, but very soon; down the line, you will start discussing strained silicon also. So, we will postpone that discussion for the time being. This uses 1.2 nanometer oxide, this is even thinner than previous generation and nickel silicide; instead of cobalt silicide, it is nickel silicide.

(Refer Slide Time: 45:26)



Now, there is a some more optimization and so, this is the cross section of a transistor; channel is strained silicon 1.2 nanometer gate oxide, 50 nanometer gate length, nickel silicide is shown here at the source drain and at the polysilicon. And they have been able to make a 60 SRAM cell using the technology.

(Refer Slide Time: 45:48)



This here just shows that you have silicon and 1.2 nanometer oxide and polysilicon gate and typical off current verses on current; for N and P channel transistors.

(Refer Slide Time: 46:03)



And they are just showing here as a function of here; different technology nodes, we already looked at 180 nanometer technology, 130 nanometer in 180 nanometer; you had 100 nanometer, 130 nanometer, 70 nanometer gate length; 90 nanometer, 50 nanometer gate length; this is what I said, now we have a departure in terms of what gate length you print is actually smaller than the so called technology definition; 90 nanometer.

But in the past; they were exactly identical, we have had discussion on this; in the one of the earlier lectures. Why nickel silicide? Because with scaling; cobalt silicide will start increasing the resistance, whereas nickel silicide; gives very low resistance and that is why they have gone to nickel silicide.

(Refer Slide Time: 46:43)



Again multi level metal interconnects; now instead of fluorine doped oxide, they have done carbon doped oxide.

(Refer Slide Time: 46:57)



I know again a new improvement; that is how you need newer and newer materials; k has come down to 2.9, instead of 3.5. Again to decrease the R C delay of your metal interconnects.

(Refer Slide Time: 47:18)



And they again show that; with carbon doped oxide, you can get significantly lower RC delay as you start decreasing the pitch. As compared to aluminum with fluorine doped oxide, copper with fluorine doped oxide and this is copper with carbon doped oxide; which is much better. And they have been able to do an SRAM here, and this SRAM operates at 2 gigahertz; which is much higher frequency.

I suspect that the previous paper that probably should have been 1.6 giga hertz; that could be a typo there; rather than 1.6 mega hertz; now that is more like it. Because on 180 nanometer, you had a 900 mega hertz, very close to gigahertz and then you have 1.6 gigahertz and now you have about 2 gigahertz or little more than 2 gigahertz.

(Refer Slide Time: 48:51)



So, this is about 90 nanometer and if you come to 65 nanometer; all that has happened really is gate length of scale down further; 35 nanometer, 8 levels of copper interconnects, interconnect hierarchy is also becoming more complex and SRAM cell is also shrinking. Continue to use nickel silicide. And let us see if there is anything; I want to highlight here.
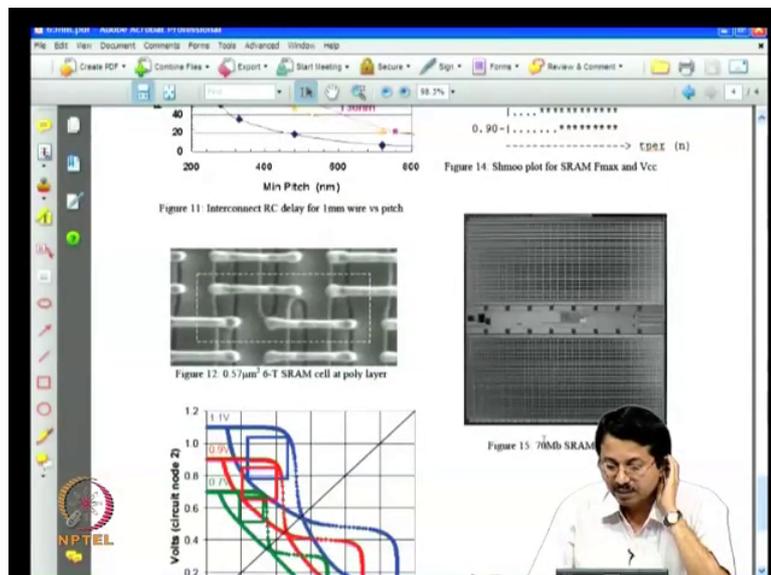
(Refer Slide Time: 48:58)



So, this is a typical trend; so, on a 65 nanometer technology, the gate lengths are 35 nanometer; this is a typical cross section of a transistor.

(Refer Slide Time: 49:14)



N and P channel transistor ID, VDS characteristics and sub threshold characteristics, this is off current versus on current; typical DC performance metric. And this is; the same for PMOS transistor, this is a very nice cross section showing different levels of metal interconnects.

(Refer Slide Time: 49:38)



And they have made an SRAM, this is 70 mega bit SRAM test vehicle. And they have been able to assure that this SRAM operates at frequency is close to 2 gigahertz and that

they have been able to demonstrate; no major departures. Then at 45 nanometer; you see a very significant departure.

(Refer Slide Time: 50:01)



And what is that? For the first time, you see high k gate dielectrics; along with metal gate. Here high k metal gate for the first time; 45 nanometer logic technology and transistors feature; 1 nanometer EOT, that is equivalent oxide thickness, but they actually use a hafnium oxide gate dielectric, which is thicker than 1 nanometer. They use a high k gate dielectric, dual band edge work function metal gates; two different metals; one for N channel, one for P channel transistors. So, what they have been able to show here is the following.

(Refer Slide Time: 50:50)



You see this is gate leakage here and this is oxide thickness; EOT, when you use SiO 2; EOT and SiO 2 are one and the same. You see, from previous this is a 45 nanometer generation, at 65 nanometer generation; they were not able to scale the gate oxide thickness. Because the leakage current; as we start a decreasing the gate oxide thickness, leakage current suddenly started increasing because of direct tunneling current. And because from 90 nanometer to 65 nanometer; they did not have a mature high k gate dielectric technology, they could not scale the oxide thickness, they kept the oxide thickness as is and then the leakage current was kept at the same time, same level.

Now, with the high k technology; the EOT came down, but the leakage did not increase; the leakage also came down. Because it is a EOT of 1 nanometer, but indeed it is a thicker; physical oxide, direct tunneling is suppressed and hence you got the best of the world. You got lower EOT and lower leakage current; that is why you want to use high k dielectric. That is illustrated very nicely here, showing the trends in different technologies. So, this is the high k and metal gate is different for NMOS and PMOS.

(Refer Slide Time: 52:28)



So, what is the typical process flow? They do a STI; Shallow Trench Isolation, wells and adjust the Vts, they deposit high k gate dielectric using Atomic Layer Deposition; ALD is Atomic Layer Deposition of the dielectric. They do a disposable gate; metal gate, remember the discussion that we had. So, polysilicon is done as usual, gate patterning, source drain extension, spacers; everything is done.

Source drain formation, nickel silicidation and then you deposit an insulator. Then poly opening you yetch-back the poly; poly is removed that is what we mean by disposable metal gate. Poly was sitting there as a place holder so that you get self aligned transistors. Then you do one metal for PMOS transistor; which has a work function appropriate with P channel transistor. Do patterning of the gate and then again do a NMOS work function metal deposition.

You need two metals for defining two different kinds of transistor and you complete the transistors and then of course, you do the multi level metal interconnects.

(Refer Slide Time: 53:40)



So, this is again to show that with high k; your gate leakage has come down very dramatically. If you had used the silicon oxide or even so called nitride silicon oxide, you would had a very leakage current, but now that is not the case. This is the typical Vt versus gate length; threshold voltage role off, that we have seen; measured that low drain voltage and high drain voltage; that is for N channel and this is for P channel transistors.

(Refer Slide Time: 54:16)



Typical off current versus on current plots; looks similar, everything looks similar.

(Refer Slide Time: 54:21)



And now typical delay versus off current; remember? This is how it looks, when you have very large off current; it means the smallest length transistor has a highest leakage current. And hence because you have build the inverter using the smallest length transistor you have the best delay that you can get. And the delays are now of the order of very low delays, 5 picoseconds or even less than 5 picoseconds.
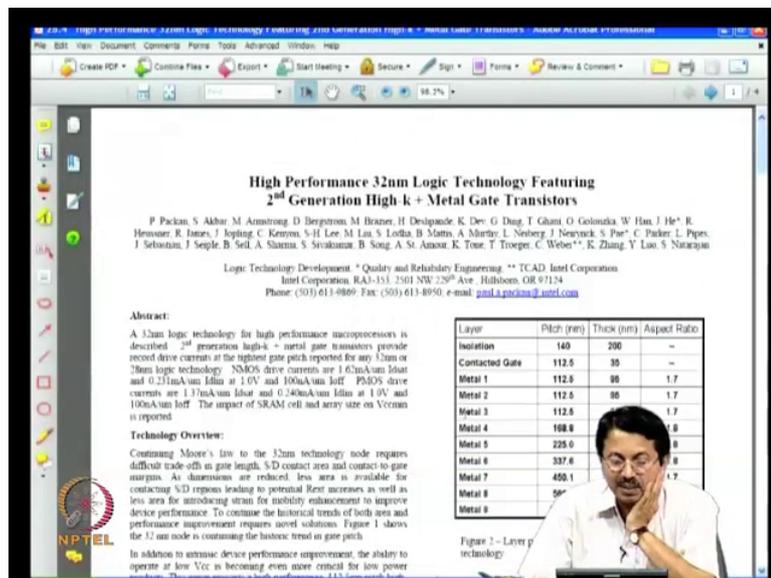
Reliability of the transistors; again they have characterized TDDB by accelerated testing; extrapolation, just to show that these high k gate dielectrics are as reliable as silicon oxide.

(Refer Slide Time: 55:07)



Then they have actually made, not only SRAM; in this case they have actually made a microprocessor. They have made a single core microprocessor and dual core microprocessor and been able to show that; you can actually get a fairly good circuit working with all these new technology; high k and metal gate so, on and so forth.

(Refer Slide Time: 55:37)



And that is what happened in 45 nanometer technology and in 32 nanometer technology; with high k and metal gate continued, it continues to be high k and metal gate transistors, we have stuck to that right now. Going forward, we will continue to use high k and metal

gate, but this called second generation of high k and they have tried to do some optimization; improve the hafnium oxide performance, improve the metal gate performance and so on and so forth.

So, 32 nanometer; the actual gate length is little lower than 32 nanometer.

(Refer Slide Time: 56:12)



This is just showing that previous generation technology and this is 32 nanometer technology.

(Refer Slide Time: 56:23)

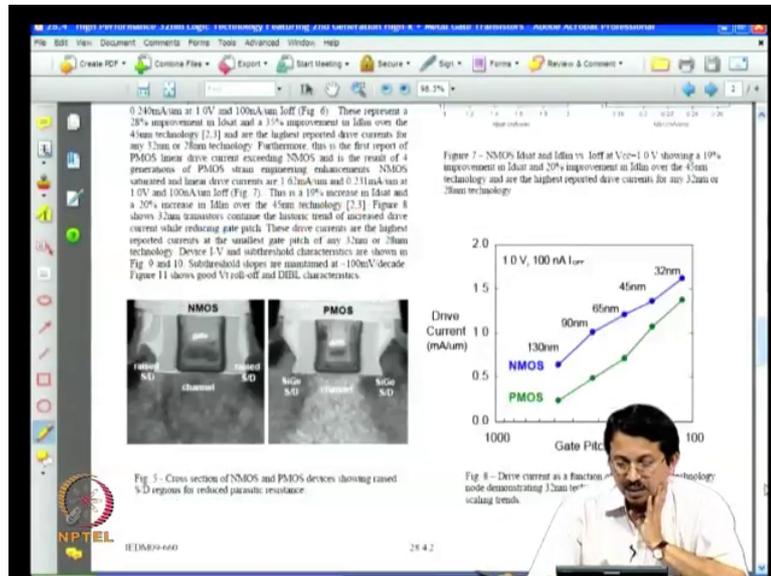And off current versus on current, the typical plots that we have seen. All these are off current versus on current plots, we trying to indicate that; they have been able to optimize these transistors.

(Refer Slide Time: 56:36)



For a given off current; they have been able to get as large on current as possible. Typical output characteristics; IDS versus VDS; notice that the difference between N and P has come down.
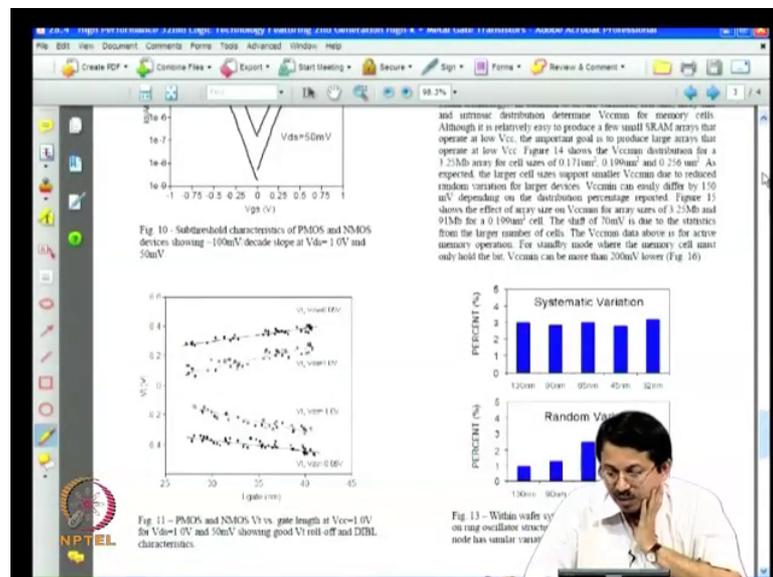
(Refer Slide Time: 56:44)

If you recall 180 nanometer technology, your P channel current was almost 50 percent of the N channel current. Now, using the so called strain engineering which; as I said we will be discussed in one of the future classes. We are able to enhance both electron and whole mobilities and so much so that the whole mobility can approach very close to electron mobility.

And that is why you do not see a huge difference between electron and hole mobility and electron and hole current; N channel and P channel current.

(Refer Slide Time: 57:32)



That difference has come down; this is a typical sub threshold plot that you would see. And typical Vt versus length plot; Vt is decreasing with decreasing channel length. So, all these is consistent with what we have been discussing and they have again made an SRAM and they demonstrate that; SRAM works fairly well with this 32 nanometer technology.

(Refer Slide Time: 57:46)



(Refer Slide Time: 57:48)
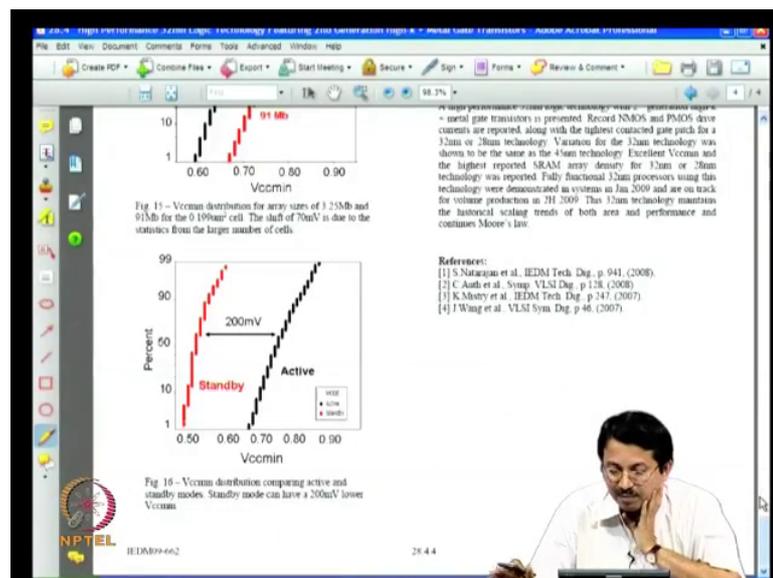


So, what we have been able to do in this lecture is really to look at representative CMOS logic technology, not memory technology; these technologies are not for DRAM fabrication. Not for flash memory fabrication, not for even SRAM fabrication, but to make a very high performance microprocessor. So, their performance is the key; you want to make the transistor, which gives a lowest delay.

And various techniques we have seen, all source drain engineering, channel engineering is consistently use in all these transistors. And as we have continued to scale the

technologies, we have transition from silicon oxide to high k gate dielectrics; to metal gates and so on and so forth.

Similarly, in the interconnect arena; we have gone from aluminum and conventional silicon oxide to dope silicon oxide, copper and there by minimize the R C delay of interconnects very significantly. So, hopefully now you should be able to read up any CMOS technology paper and be able to understand that; based on all the discussion that we have had so far.

So, we will stop this lecture today and starting from next lecture, we will actually start looking at electrical characterization; that is C V characterization and I V characterization.