

Semiconductor Devices and Circuits
Prof. Sanjiv Sambandan
Department of Instrumentation and Applied Physics
Indian Institute of Science, Bangalore

Lecture – 42
Scaling of MOSFETs

(Refer Slide Time: 00:15)

Scaling of MOSFETs

To reduce the size of the integrated circuit chips, reduce power consumption and increase speed, the MOSFETs are scaled down.

1970 - 2250 transistors in Intel 4004
 2015 - 10000000000 transistors in SPARC M7
 1970 - 10 microns channel length
 2015 - 40nm channel length

Scaling implies: Reduce L , W and t_{ox} by a factor of α .

How do we scale down?
 There are two methods of scaling:

(i) Constant field scaling:
 Need to keep field constant. If we scale down L , W and t_{ox} , we should scale down applied voltages. Voltage levels not compatible with previous versions.

(ii) Constant voltage scaling
 Power supply voltage not reduced, therefore compatible with previous versions. But field becomes very large in the device.

So, now let us start off, on a new topic, which is got to do the Scaling of the MOSFETs, and this is a particular interest as it is the symbol, or you know, it is very symbolic of the advancement of your semiconductor technology. So, it is quite, famous particularly in the form of Moore's law. Although surprisingly, I am never referred to that here, which basically, says which defines a pathway or a roadmap to scaling, which says that you know the number of transistors in any system or any chip, doubles every so many years; but just to give you an example.

So, what is scaling essentially before we start talking about scaling, what does what does scaling mean. So, we have a MOSFET here which is, you know drawn in a very simple manner, that is your source, that is the drain and that is the gate and that is the insulator, that is got some thickness t_{ox} . There is a channel length between source and drain and the MOSFET has got some channel width.

So, what do you mean by scaling? Scaling simply means, reducing the dimensions so, what we want to do is, when we say we are going to scale the MOSFET down by a factor of alpha, what we mean is, we are going to reduce the channel length by a factor of alpha. We are going to reduce the channel width by a factor of alpha and we are going to reduce the thickness by a factor of alpha.

So, that is what scaling implies and why do we need to do scaling; so, that you can make the transistor smaller and so that, you can have more transistors in a given layout area and by scaling down the transistor you improved, the speed and the performance and I mean, the speed and performance in the power. We change the power consumption, depending on how you scale, of course, and we will look at that and the most importantly you can fit in a larger number of transistors in a given layout area.

So, there are lots of benefits. These kind of benefits to scaling and that is what is led to you know, the microprocessors, in the in the computers becoming faster and, you know, it is led to, mobile technology etcetera etcetera. So, it is quite an important feature of any, technology has to how rapidly and you know how effectively you can scale down to improve performance.

So, just to give you an example, I hope, I am right on all these, but, if I am not, I should, I would probably correct it, but in the 1970s you know. So, Intel 4004 processor had about 2000 transistors in it ok. And in 2015, that is what, 45 years ok. So, in half centuries time, the number of transistors, has increased to these many.

That is a massive increase and in the 1970s going back to the 1970 the, there is a number of transistors. The transistors had a typical channel length of 10 microns and in 2015, there the channel length of the transistors about 40 nanometers ok. So, that is going, that is going by 3 orders of magnitude lower in terms of channel lengths. So, these kinds of advancements are possible by scaling down MOSFETs, in a reliable manner.

So, what we are going to do here is, we are going to look at, what scaling is introduce it. Talk about two different kinds of scaling strategies and what its effect is, ok. Some of the key features and the parasitic that rise along with scaling that come along with scaling etcetera.

(Refer Slide Time: 04:19)

Scaling of MOSFETs

To reduce the size of the integrated circuit chips, reduce power consumption and increase speed, the MOSFETs are scaled down.

1970 - 2250 transistors in Intel 4004
2015 - 10000000000 transistors in SPARC M7
1970 - 10 microns channel length
2015 - 40nm channel length

Scaling implies: Reduce L , W and t_{ox} by a factor of k .

How do we scale down?
There are two methods of scaling:

(i) Constant field scaling:
Need to keep field constant. If we scale down L , W and t_{ox} we should scale down applied voltages. Voltage levels not compatible with previous versions.

(ii) Constant voltage scaling
Power supply voltage not reduced, therefore compatible with previous versions. But field becomes very large in the device.

Mour'idan

W

V_{gs}

alpha

So firstly, how do we scale down ok. So, as I told you, you reduce L W and t_{ox} by a factor of α ok. It does not matter. So, there are two strategies, one is something called as constant field scaling and the other is something called as constant voltage scaling and of course, you can have a hybrid of these two, but we will just introduce these two strategies. In the case of constant fields scaling the electric fields, inside the device are kept constant ok, and what does that mean? So, let us say, we have scaled down, we have scaled down the MOSFET from L to L by α .

So, there were there was a certain electric field in this MOSFET ok. So, there was a there was a gate field, there is a source to drain field etcetera and we want to keep the magnitude or these field strengths the same even though the transistors have been made smaller. And how do we achieve that? We achieve that by simultaneously scaling down the voltages ok. So, you want to keep the fields the same and therefore, the voltages have to go down along with your dimensions.

(Refer Slide Time: 05:39)

Scaling of MOSFETs

To reduce the size of the integrated circuit chips, reduce power consumption and increase speed, the MOSFETs are scaled down.

1970 - 2250 transistors in Intel 4004
 2015 - 1000000000 transistors in SPARC M7
 1970 - 10 microns channel length
 2015 - 40nm channel length

Scaling implies: Reduce L , W and t_{ox} by a factor of k .

How do we scale down?
 There are two methods of scaling:

(i) Constant field scaling:
 Need to keep field constant. If we scale down L , W and t_{ox} , we should scale down applied voltages. Voltage levels not compatible with previous versions.

(ii) Constant voltage scaling:
 Power supply voltage not reduced, therefore compatible with previous versions. But field becomes very large in the device.

The other idea is something called as constant voltage scale which, in this case even despite your scaling down of lengths the voltages remain constant. In order for you, to continue having compatibility with regards to the power supply required for operation etcetera and; obviously, in this case, the electric fields are going to increase as you scale down the MOSFETs.

(Refer Slide Time: 05:57)

Scaling of Transistors

Parameter	Constant Field	Constant Voltage
Channel Width, W	$1/k$	$1/k$
Channel Length, L	$1/k$	$1/k$
Oxide Thickness, t_{ox}	$1/k$	$1/k$
Area	$1/k^2$	$1/k^2$
Electric field	1	k
Voltage	$1/k$	1
Insulator capacitance/area, C_{ox}	k	k
Total Gate Capacitance	$1/k$	$1/k$
Current	$1/k$	k
Power	$1/k^2$	k

$k=4$
 $L=8\mu m$
 $L'=2\mu m$
 $W' = \frac{W}{4}$
 $t_{ox}' = \frac{t_{ox}}{4}$

So, you have a quick table here. So, I will just be prepared easy table, for see as to what happens when you use constant fields scaling and constant voltage scaling. So, here we have

scaling down by a factor of k . So, let us say, we say k is equal to 4, what that means is, my channel and let us say, my initial channel length was 8 microns, what that means is, if I scale down by a factor of 4, my new channel length has become 2 microns ok, and my new channel width has become the old channel width by 4. The new t_{ox} has become the old t_{ox} by 4 etcetera-etcetera. So, that is what scaling down by a factor of k implies.

(Refer Slide Time: 06:49)

Scaling of Transistors

Parameter	Constant Field	Constant Voltage
Channel Width, W	$1/k$	$1/k$
Channel Length, L	$1/k$	$1/k$
Oxide Thickness, t_{ox}	$1/k$	$1/k$
Area	$1/k^2$	$1/k^2$
Electric field	1	k
Voltage	$1/k$	1
Insulator capacitance/area, C_{ox}	k	k
Total Gate Capacitance	$1/k$	$1/k$
Current	$1/k$	k
Power	$1/k^2$	k



Now, how does it affect the simple parameters? Channel width of course, goes down by a factor of k , whether its constant voltages are constant field. Channel length goes down by a factor of k , with its constant voltage a constant field, oxide thickness goes down by a factor of k , with a channel length, whether it is constant field a constant voltage.

The area is basically, your W into L essentially of course, there are lot more things to the area, because you have to have your pads, you have to have your interconnects, but to first order the size of the transistor. If you are familiar with transistor layout, then you have a certain channel length, you have a certain channel width and therefore, the area simply scales as W over L plus of course, there are many other, aspects to the area.

So that, we will go down as 1 by k square since both my W and L reduced by a factor of k ; so, that is area and the electric field in constant field scaling. The electric field is not supposed to change, that is the whole idea. So, that does not scale and remains constant at 1 and if you want to keep the electric field constant, you have to reduce your voltage and the

electric field is essentially, the voltages divided by all the lengths, and since the lengths have scaled down by a factor of k, we also need to scale down the voltages by a factor of k.

(Refer Slide Time: 07:49)

Scaling of Transistors

Parameter	Constant Field	Constant Voltage
Channel Width, W	1/k	1/k
Channel Length, L	1/k	1/k
Oxide Thickness, t_{ox}	1/k	1/k
Area	1/k ²	1/k ²
Electric field	1	k
Voltage	1/k	1
Insulator capacitance/area, C_{ox}	k	k
Total Gate Capacitance	1/k	1/k
Current	1/k	k
Power	1/k ²	k

Handwritten notes:
 $f = \frac{V/k}{L/k} = \frac{V}{L}$
 $C_{ox} = \frac{\epsilon_{ox}}{t_{ox}}$
 $C_{ox} = \frac{\epsilon_{ox}}{t_{ox}} = \frac{\epsilon_{ox}}{t_{ox}/k} = k \frac{\epsilon_{ox}}{t_{ox}}$
 Diagrams of a transistor showing dimensions W, L, and t_{ox} with arrows indicating scaling directions.

So, the voltage in the case of constant field scaling will scale down by a factor of k. On the other hand, the voltage in a constant voltage scaling case will not change. It has to be constant and therefore, the electric fields all increase by a factor of k.

(Refer Slide Time: 08:33)

Scaling of Transistors

Parameter	Constant Field	Constant Voltage
Channel Width, W	1/k	1/k
Channel Length, L	1/k	1/k
Oxide Thickness, t_{ox}	1/k	1/k
Area	1/k ²	1/k ²
Electric field	1	k
Voltage	1/k	1
Insulator capacitance/area, C_{ox}	k	k
Total Gate Capacitance	1/k	1/k
Current	1/k	k
Power	1/k ²	k

Handwritten notes:
 $C_{ox} = \frac{\epsilon_{ox}}{t_{ox}}$
 $C_{ox} = \frac{\epsilon_{ox}}{t_{ox}} = \frac{\epsilon_{ox}}{t_{ox}/k} = k \frac{\epsilon_{ox}}{t_{ox}}$
 $f = \frac{V/k}{L/k} = \frac{V}{L}$
 Diagrams of a transistor showing dimensions W, L, and t_{ox} with arrows indicating scaling directions.

So, what about the insulator capacitance per unit area; so, that is C_{ox} that is the capacitance per unit area is nothing, but ϵ_{ox} by t_{ox} and our t_{ox} scaled down by a factor of k and

therefore, the insulator capacitance per unit area has to increase by a factor of k and that is true, whether this constant field or constant voltage scaling.

Now, what about the total gate capacitance and what do you mean by the total gate capacitance? So, this is the total this is the insulated capacitance per unit area and therefore, the gate capacitance is the total capacitance; so, C_{ox} into $W L$ which is the area of the device. So, since C_{ox} increased by a factor of k and W decreased by a factor of k and L decreased by a factor of k irrespective, whether it's constant voltage or constant field scaling the total gate capacitance reduces by a factor of k, what about the current in the devices.

(Refer Slide Time: 09:43)

Scaling of Transistors

Parameter	Constant Field	Constant Voltage
Channel Width, W	$1/k$	$1/k$
Channel Length, L	$1/k$	$1/k$
Oxide Thickness, t_{ox}	$1/k$	$1/k$
Area	$1/k^2$	$1/k^2$
Electric field	1	k
Voltage	$1/k$	1
Insulator capacitance/area, C_{ox}	k	k
Total Gate Capacitance	$1/k$	$1/k$
Current	$1/k$	k
Power	$1/k^2$	k

Handwritten notes:
 $I_{lin} = \mu C_{ox} \frac{W}{L} (V_{gs} - V_t - \frac{V_{ds}}{2}) V_{ds}$
 Diagram of a transistor with dimensions W , L , and t_{ox} .
 Area WL

Now, whether it is a linear mode or saturation mode the current has a square dependence on the voltage. So, if it is a linear you say its $\mu C_{ox} W$ over L . So, let me write this properly; so, in case of linear, we say that it is $\mu C_{ox} W$ over L , into V_{gs} minus V_t minus V_{ds} by 2 into V_{ds} . So, that is the current ok. So, you can see that there is a voltage into voltage term. So, that is a V square term V square, the voltage square dependence and in saturation of course, its $\mu C_{ox} W$ over L by 2, into V_{gs} minus V_t the whole square.

(Refer Slide Time: 10:15)

Scaling of Transistors

Parameter	Constant Field	Constant Voltage
Channel Width, W	1/k	1/k
Channel Length, L	1/k	1/k
Oxide Thickness, t_{ox}	1/k	1/k
Area	1/k ²	1/k ²
Electric field	1	k
Voltage	1/k	1
Insulator capacitance/area, C_{ox}	k	k
Total Gate Capacitance	1/k	1/k
Current	1/k	k
Power	1/k ²	k

Handwritten notes on the left: $I_{lin} = \mu C_{ox} \frac{W}{L} (V_{gs} - V_{th})^2$ and $I_{sat} = \mu C_{ox} \frac{W}{2L} (V_{gs} - V_{th})^2$.
 Handwritten notes on the right: A diagram of a transistor with width W and length L, and the text "Area WL".

So, again it is a voltage square dependence ok. So, what happens if you scale down? So, what are the parameters here that scale, the mobility not so much W and L scale, but W divided by L does not scale. That is the aspect ratio since W scales by k and L scales down by k. W over L does not scale C_{ox} does scale. So, C_{ox} how does, C_{ox} scale, in the case of constant field scaling, it C_{ox} increases by a factor of k.

(Refer Slide Time: 11:07)

Scaling of Transistors

Parameter	Constant Field	Constant Voltage
Channel Width, W	1/k	1/k
Channel Length, L	1/k	1/k
Oxide Thickness, t_{ox}	1/k	1/k
Area	1/k ²	1/k ²
Electric field	1	k
Voltage	1/k	1
Insulator capacitance/area, C_{ox}	k	k
Total Gate Capacitance	1/k	1/k
Current	1/k	k
Power	1/k ²	k

Handwritten notes on the left: $I_{lin} = (k C_{ox}) \frac{W}{L} (V_{gs} - V_{th})^2$ and $I_{sat} = k C_{ox} \frac{W}{2L} (V_{gs} - V_{th})^2$.
 Handwritten notes on the right: A diagram of a transistor with width W and length L, and the text "Area WL".

So, C_{ox} becomes k time C_{ox} and how does voltage change, with in case of constant field scaling, the voltage reduces has to reduce by a factor of k and therefore, in the case of

constant field scaling, the voltage reduces by a factor of k and therefore, in the case of in the case of constant field scaling. So, whether it is here or here it is the same. So, in the case of constant field scaling the current has to reduce by a factor of k .

(Refer Slide Time: 11:59)

Scaling of Transistors

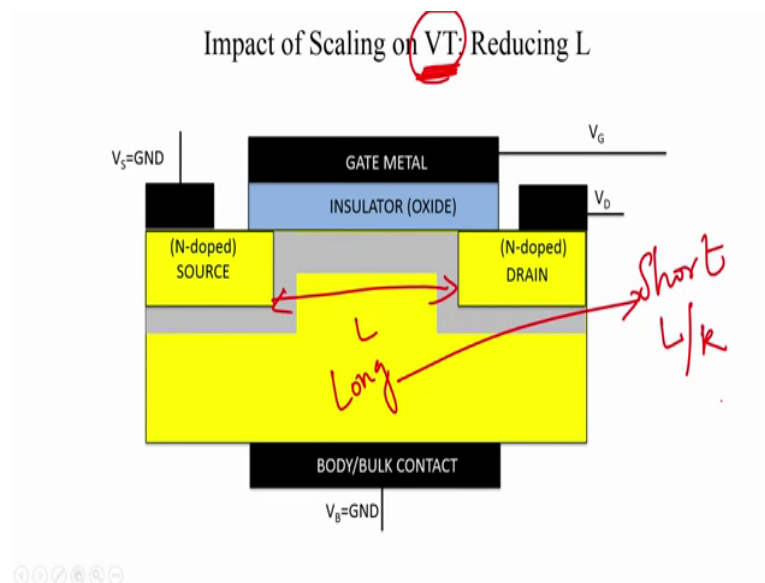
Parameter	Constant Field	Constant Voltage
Channel Width, W	$1/k$	$1/k$
Channel Length, L	$1/k$	$1/k$
Oxide Thickness, t_{ox}	$1/k$	$1/k$
Area	$1/k^2$	$1/k^2$
Electric field	1	k
Voltage	$1/k$	1
Insulator capacitance/area, C_{ox}	k	k
Total Gate Capacitance	$1/k$	$1/k$
Current	$1/k$	k
Power	$1/k^2$	k

Handwritten notes:
 Left side: $I \sim (C_{ox} V^2) / L$, $P = I \cdot V$
 Right side: Diagram of a transistor with dimensions W , L , and t_{ox} .
 Bottom right: $A_{area} \sim W \cdot L$

So, you can see its V square by k square into $k C_{ox}$ and therefore, the current reduces by a factor of k on the other hand, in the case of constant voltage scaling. So, let us just write, I just keep things simple. We just write the current, it is current scale, scalable the scaling aspects of the current is dependent on C_{ox} into V square we, because the rest of the parameters do not change. So, in the case of constant voltage scaling V does not scale.

So, that remains 1 whereas, C_{ox} has increased by a factor of k . So therefore, the current increases by a factor of k ok. So, you can see that, in a constant field scaling case, the current in the transistor decreases by a factor of k , in the constant voltage scaling case. The current in the transistor increases by a factor of k and what about power consumption power is nothing, but your current into voltage IV . And, since all the voltage is reduced by a factor of k , in constant field scaling and the current reduced by a factor of k , the power reduces by a factor of k square. Whereas, in the case of constant voltage scaling, the current increased by a factor of k , the voltage remained invariant and therefore, the product scales up by a factor of k , ok.

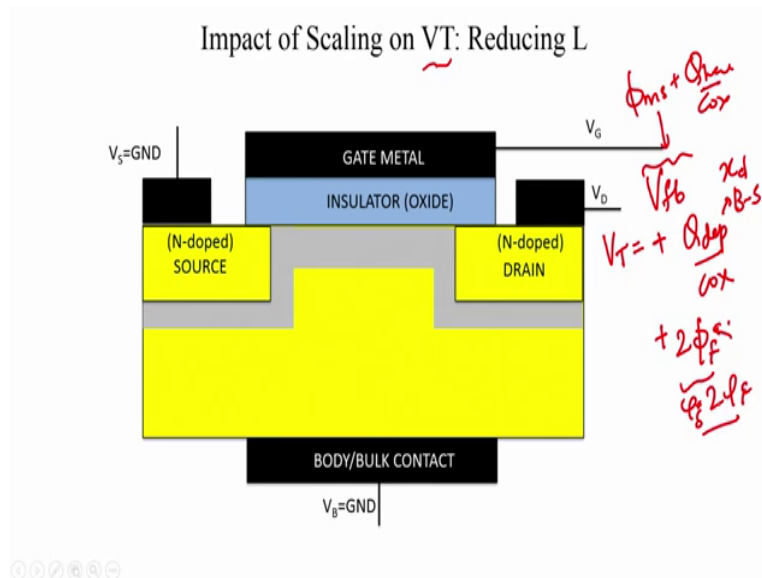
(Refer Slide Time: 11:33)



So, here the power increases, here the power comes down ok, here the current increases here the current comes down. So, this is how you handle scaling relations of course, there are many other parameters, one can consider you know delay etcetera-etcetera, but we have just talked about some other basic or the key parameters. So, what is the impact of scaling ok; so, the first impact of scaling down is got to do with the threshold voltage ok. So, we will talk about this first. So, is the threshold voltage influenced by scaling of the MOSFET ok.

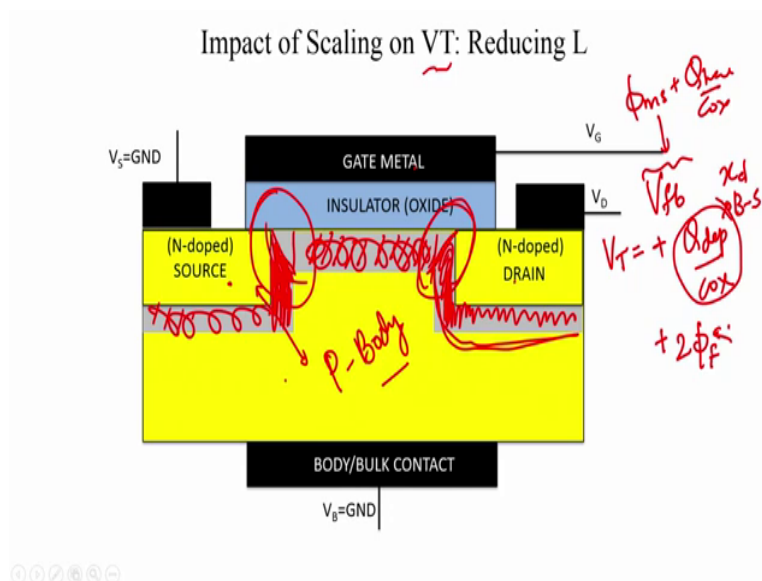
So, what we are going to do is, we are going to start with a regular long channel device. So, it is a long channel device and then we are going to head, towards a short channel device, which is the scaled down version short channel device, and we are going to see what has to, what happens to the insulator. So, before we head that, how could the threshold; we are going to see what happens to the threshold voltage. Sorry I am probably getting 2, to use to saying these terms and I am, making mistakes, but nevertheless, let us just carry on.

(Refer Slide Time: 14:39)



So, we have the threshold voltage and what is the threshold voltage depend on. So, once again let me repeat, it has got these three milestones. You need to first get your flat band voltage which you now know that, it depends on ϕ_{ms} plus any trapped charge in the insulator. Then you have to have your Q_{dep} by C_{ox} ok, and this depends on the varying x_d , it depends on, the body to source voltage etcetera and then you need to achieve your $2\phi_s$ ok, but here again it depends on the body to source voltage, which is your ϕ_s has to be equal to $2\phi_s$. So, this is when you hit threshold voltage ok.

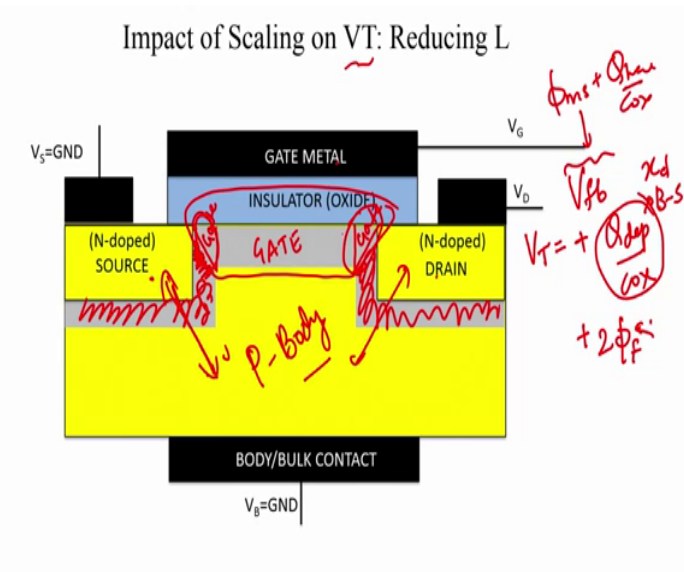
(Refer Slide Time: 15:31)



Now, let us look at this depletion term. So, what exactly are we depleting in a MOSFET. So, here you already have a NP-junction. So, this is a p-type body; So, you have a n-type source a n-type drain and you have a P-type body and you already have a depletion layer here, because of this diode.

So, this region is already depleted for you. The gate has to deplete this region of the semiconductor ok, what about this region here, that region is shared between the gate and the PN-junction. So, this PN-junction at equilibrium will help the gate and depleting that. So, let us say, half of that region is depleted by the PN-junction and the remaining half is the job of the gate. So, the gate gets some help in that region, in these regions. It is only the PN-junction and that region is only the gate and in this region the PN-junction helps the gate depleted, but in a long channel device ok.

(Refer Slide Time: 16:45)

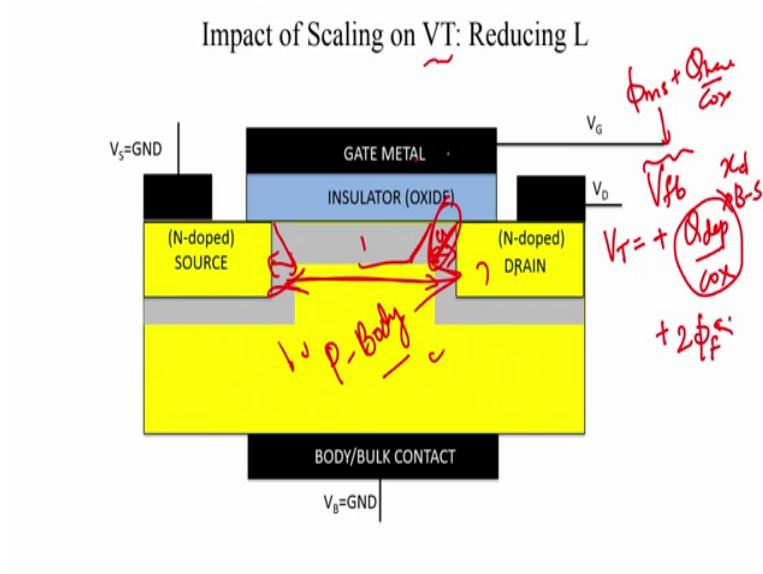


So, let me just write this down neatly, in case my; so, this is the job of the gate. The gate needs to deplete this region. Half of this region is also the job of the gate, because this region is shared between the PN-junction and the gate and everything else the remainder half and all these regions are all depleted by this PN-junction between source and drain and bulk. So, all these are jobs of the PN-junction.

So, the gate essentially, all the gate charges go in to deplete that particular region. Now, if it is a long channel device, these regions here are as negligible portion of the entire area that gate

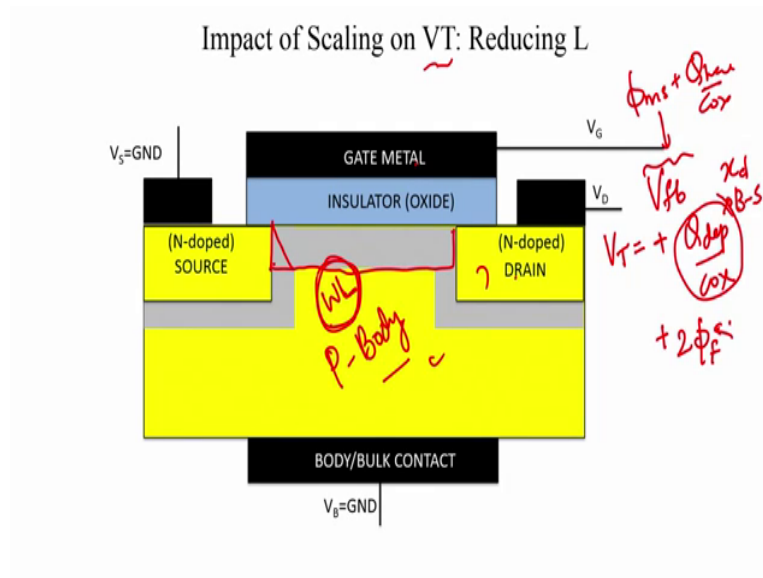
has to deplete. So yes, if you have a very long channel device and you reduce the channel length by 5 percent, the gate is not going to see too much of, in effect ok.

(Refer Slide Time: 17:57)



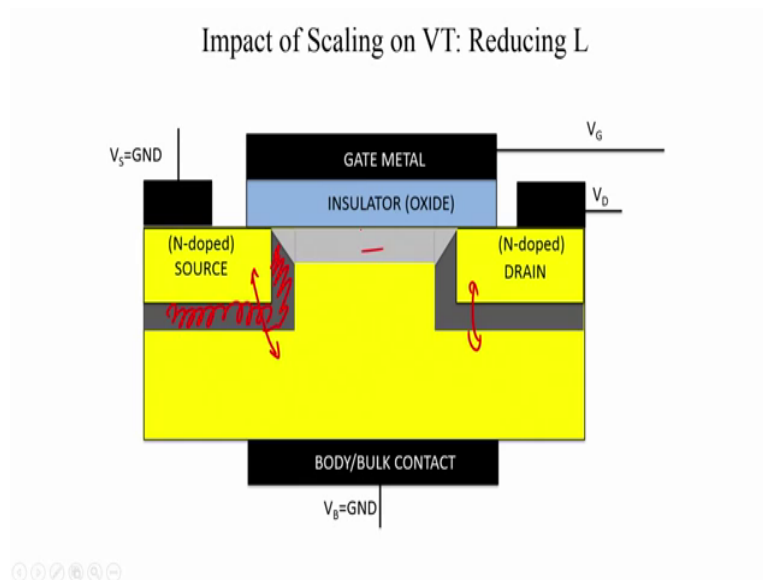
So, what do you mean by long channel device? Firstly, long channel device implies that, the distance between these channels ok. The channel lengths, sorry the distance between these two electrons, the channel length is much larger as compared to all these depletion regions, much larger than these ok. And in that case, it really does not make too much of a difference, whether the gate depleted that much or whether depleted this much, because this portion is a small percentage of the entire channel area ok.

(Refer Slide Time: 18:33)



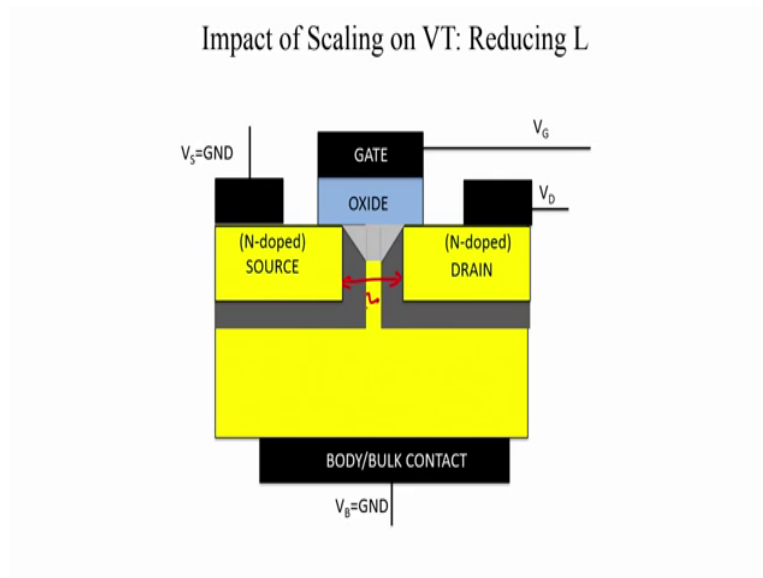
So, the gate can be said to be depleting, the entire area, which is about W into L ok. So, that is the area, the gate has to deplete.

(Refer Slide Time: 18:49)



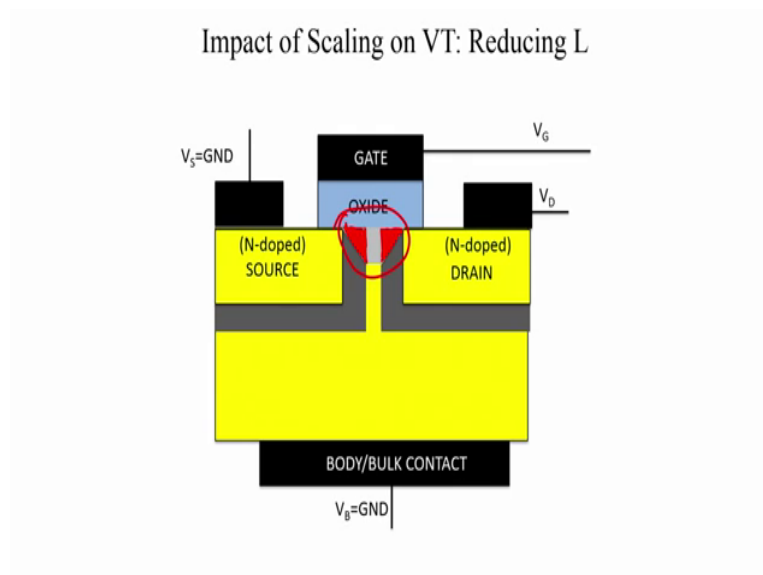
Now, what happens; if start reducing the device, if you start reducing the device ok. So, here I had a nice picture and I did not use it. This dark region shows, the portion of the, portion depleted by the PN-junction, whereas, the gray region shows that, portion depleted by the gate.

(Refer Slide Time: 19:11)

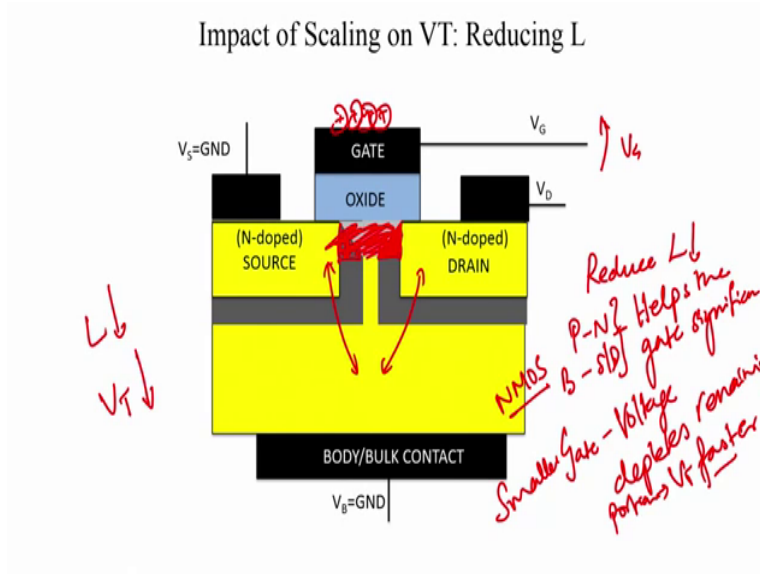


Now, if you start scaling down the channel length, so now, the channel length is reduced quite a bit ok. So now, this depletion regions are not that small, we cannot really ignore them. They have gotten very close to each other.

(Refer Slide Time: 19:27)



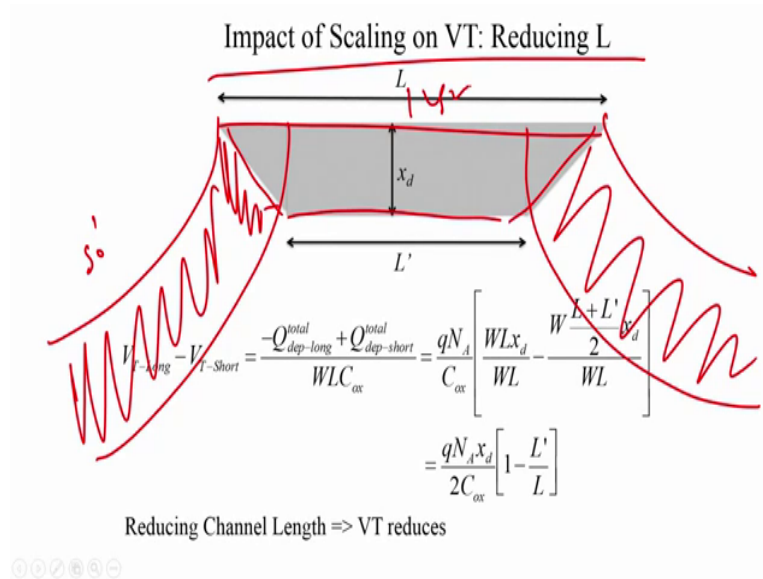
(Refer Slide Time: 19:53)



So, now this region suddenly forms a significant part of the total depletion charge, right, and now, we need to start taking the count. This area, quite carefully in order to calculate, all are make calculations on the threshold voltage, and since the source this PN-junction now, takes up a larger percentage or the larger share of the workload, in trying to deplete the gate has to do a smaller job and therefore, the gate depletes that region faster at lower V_g s, as compared to the long channel device. So, in other words, as we reduce the channel length, the PN-junctions between that is between the source drain and body the PN-junctions in case. In fact, the other way, it is the body and source or drain ok. Since, we are talking about NMOS we are always discussing the n MOSFET we have never discussed the p MOSFET.

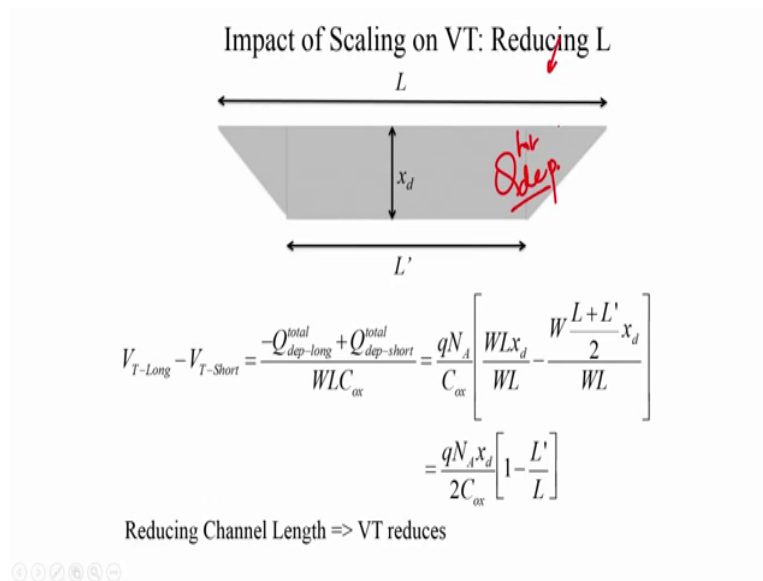
So, far; The PN-junction helps the gate significantly ok. So, that is the key significantly, and it depletes, it performs half of the workload it depletes it helps the gate deplete a lot of the semiconductor. And therefore, the gate can a smaller gate voltage a smaller gate voltage depletes the remained a portion of the semiconductor and takes the semiconductor to threshold and threshold voltage is reached faster. So, in other words, you need to add a lot less charge on the gate, in order to completely deplete this entire engine and head towards inversion ok. So, therefore, as the channel length reduces the threshold voltage reduces and if you want to see this mathematically, we just take the small analysis here.

(Refer Slide Time: 21:47)



So, let us just take the depletion region of the gate the region of the semiconductor depleted by the gate. So, that is where your source drain would be and that is how the depletion region would be around. Your source drain and this part of it is all taking care. It is all depleted by the source and drain and it is only this region. Here, that is depleted by there gate and that is the insulator ok. So, that is the picture.

(Refer Slide Time: 22:31)



Now, what is this area? You know, how much of this volume, this gate depleting what is the total depletion charge. What is the total amount of charge? I must add to the gate, in order for

me to get a total depletion charge over here and by total I mean in coulombs. So, we are not talking about per unit area or per unit volume total employs in coulombs ok.

(Refer Slide Time: 22:47)

Impact of Scaling on VT: Reducing L

$$V_{T-Long} - V_{T-Short} = \frac{-Q_{dep-long}^{total} + Q_{dep-short}^{total}}{WLC_{ox}} = \frac{qN_A}{C_{ox}} \left[\frac{WLx_d}{WL} - \frac{W \frac{L+L'}{2} x_d}{WL} \right]$$

$$= \frac{qN_A x_d}{2C_{ox}} \left[1 - \frac{L'}{L} \right]$$

Reducing Channel Length => VT reduces

(Refer Slide time: 22:59)

Impact of Scaling on VT: Reducing L

$$V_{T-Long} - V_{T-Short} = \frac{-Q_{dep-long}^{total} + Q_{dep-short}^{total}}{WLC_{ox}} = \frac{qN_A}{C_{ox}} \left[\frac{WLx_d}{WL} - \frac{W \frac{L+L'}{2} x_d}{WL} \right]$$

$$= \frac{qN_A x_d}{2C_{ox}} \left[1 - \frac{L'}{L} \right]$$

Reducing Channel Length => VT reduces

So, not per area or volume not this; so, what is the total charge? So, let us take a difference let us say. Now, what is this area? You know, how much of this volume, this gate depleting what is the total depletion charge, what is the threshold voltage of the long channel device minus the threshold voltage of the short channel device, in the case of a long channel device ok. This, it is almost the total charge is almost it is going to be $q N_A x_d$, that is a charge per unit

area into the total area which is all L has not changed too much. It is only a small percentage of L, that has changed. We can also be more accurate. Let us also do that, but let us just finish this first.

In the case of a short channel device these differences, if have the trapezium here and with one side as being L dash, the other side is being L. Now, L dash is significantly less than L whereas, in a long channel device L dash was almost the same as L. So, that is the key difference.

(Refer Slide Time: 23:55)

Impact of Scaling on VT: Reducing L

$$V_{T-Long} - V_{T-Short} = \frac{-Q_{dep-long}^{total} + Q_{dep-short}^{total}}{WLC_{ox}} = \frac{qN_A}{C_{ox}} \left[\frac{WLx_d}{WL} - \frac{W \frac{L+L'}{2} x_d}{WL} \right]$$

$$= \frac{qN_A x_d}{2C_{ox}} \left[1 - \frac{L'}{L} \right]$$

Reducing Channel Length => VT reduces

Handwritten notes: $qN_A x_d W(L)$, $L' < L$, $L' \approx L$

(Refer Slide Time: 24:07)

Impact of Scaling on VT: Reducing L

$$V_{T-Long} - V_{T-Short} = \frac{-Q_{dep-long}^{total} + Q_{dep-short}^{total}}{WLC_{ox}} = \frac{qN_A}{C_{ox}} \left[\frac{WLx_d}{WL} - \frac{W \frac{L+L'}{2} x_d}{WL} \right]$$

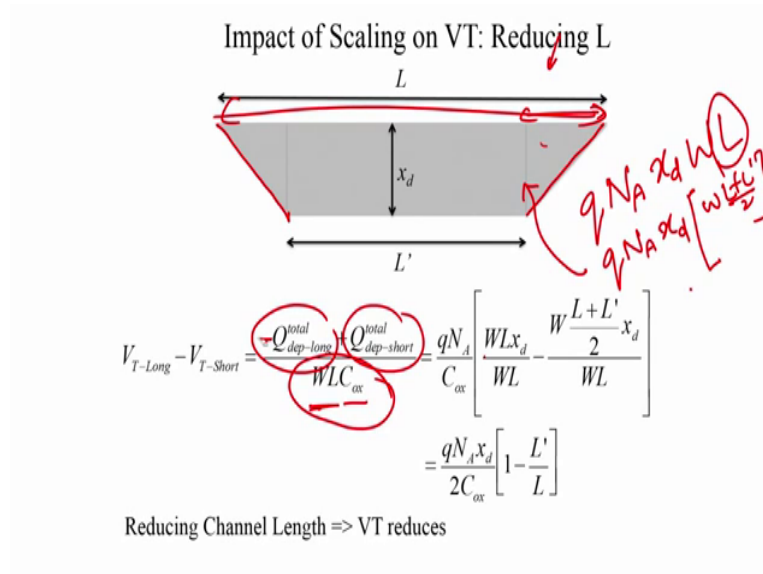
$$= \frac{qN_A x_d}{2C_{ox}} \left[1 - \frac{L'}{L} \right]$$

Reducing Channel Length => VT reduces

Handwritten notes: $qN_A x_d W(L)$, $qN_A x_d [width]$

Now, since L_{dash} is much less than L , because a larger percentage of this region is depleted by the PN-junction or this region forms a larger percentage of the, total length. Here, you have the area to be considered as the area of the trapezium. So, that $q N_A x_d$ into $W L$ plus L_{dash} by 2, in that is the area of the trapezium, divided by, I mean L plus L_{dash} by 2 is basically, your total charge, in the depletion region in the short channel device.

(Refer Slide Time: 24:51)



So, what is the difference between the V_t of a long channel device and a short channel device. It is a difference, in this depletion charge minus this depletion charge divided by the total capacitance. So, its capacitors per unit area into WL ok, and this minus sign is there, because of this minus sign is present, because of the all the charge is being negative.

(Refer Slide Time: 25:07)

Impact of Scaling on VT: Reducing L

$$V_{T-Long} - V_{T-Short} = \frac{-Q_{dep-long} + Q_{dep-short}}{WLC_{ox}} = \frac{qN_A}{C_{ox}} \left[\frac{WLx_d}{WL} - \frac{W \frac{L+L'}{2} x_d}{WL} \right]$$

$$= \frac{qN_A x_d}{2C_{ox}} \left[1 - \frac{L'}{L} \right]$$

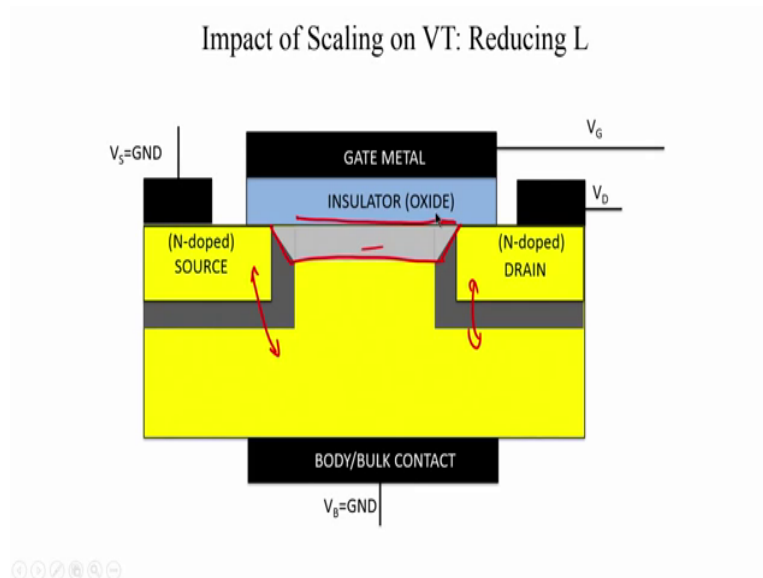
Reducing Channel Length \Rightarrow VT reduces > 0

$L' < L$

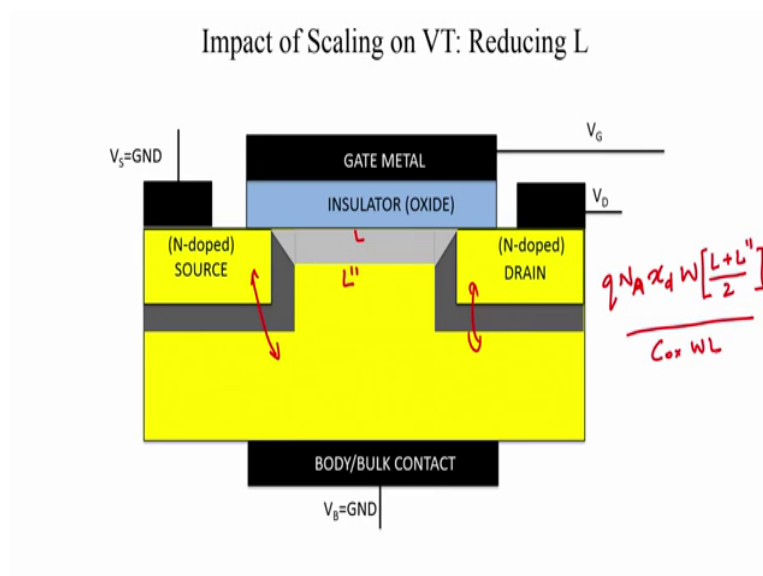
So, this total depletion charge is $q N A$ by C_{ox} into $WL \times d$ divided by WL . So, that is this term and q total depletion in the short channel device divided by the total gate capacitance is given by this term. So, essentially this reduces to an expression of this kind and if L' is less than L , this term is greater than 0, which means the V_t of the long channel device is greater than the V_t of the short channel device or in other words, if you reduce the channel length the V_t reduces ok.

Now, if you want to be more accurate, let us just go through that exercise, as well although it is not written here, we will just work it out. Let us say that, you want to say, you want to consider this trapezium even in a long channel device, that is true, because although here, the trapezium is more pronounced. It also does exist in a long channel device. It is just that, these regions are a much smaller percentage of the total length total channel length and that is the primary difference right.

(Refer Slide Time: 26:15)

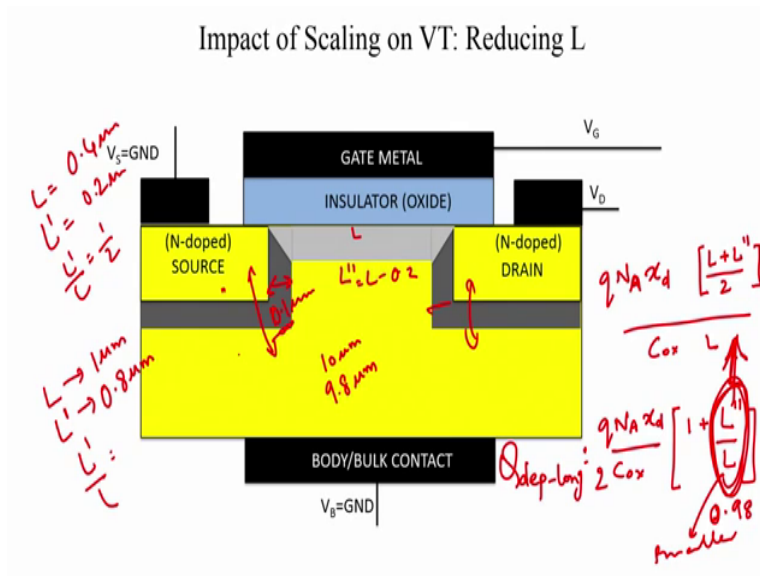


(Refer Slide Time: 26:33)



So, what is let us call this as L double dash and capital L , and what is the depletion charge. In this device it will basically, be your $q N A I$ am just taking the magnitude $q N A \times d$ into W into the area of this trapezium it should be L plus L double dash by 2, and that divided by the total gate capacitance is essentially going to be. Let us get rid of the W its essentially going to give you $q N A \times d$ by C_{ox} into 1 plus, let us say $2 C_{ox} 1$ plus 1 double dash by L .

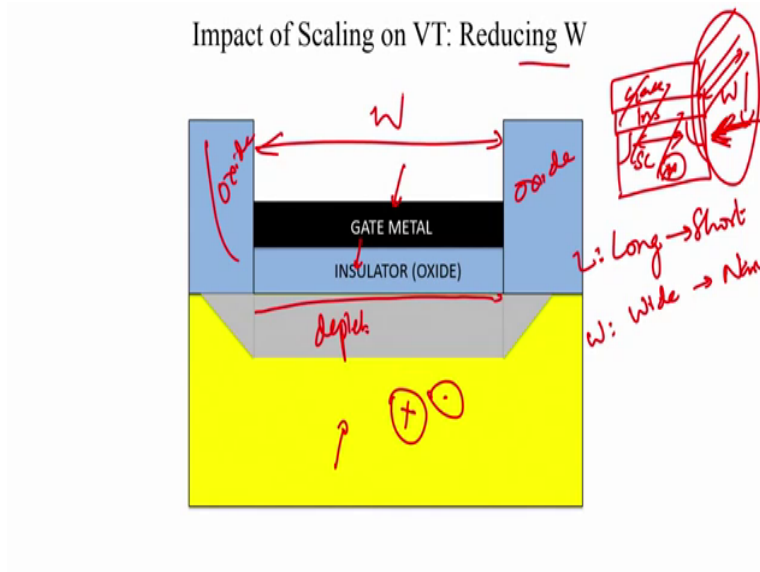
(Refer Slide Time: 27:15)



So, that is the q depletion of the long channel device in an accurate manner accurately ok. So, L say, L is 10 microns and this is say of the order of 1 e or 0.1 microns and therefore, L double dash is 9.98 microns. So, this number here is basically, 0.98. Whereas, if my L was much smaller ok. So, if on the other hand, if my L approaches, say 1 micron and this still remains to be 0.1 micron. Because, this doping concentration has not changed and my L dash for a short channel device approaches 1 minus 0.2 which is, 0.8 micron. Then my L dash by L is a smaller number and therefore, it is going to have a larger impact and further if we further reduce L .

So, let us say, we make L as 0.4 microns, as we are beginning to approach the depletion thicknesses, then my L dash should now become 0.2 microns. So, it will be L minus 0.2, because each of these regions is 0.1 and therefore, my L dash by L has now become half, so much more significant. So, you can see, that as L decreases the impact of this term begins to show up to a greater extent ok, and therefore, the threshold voltage starts becoming smaller and smaller as my L decreases.

(Refer Slide Time: 29:51)

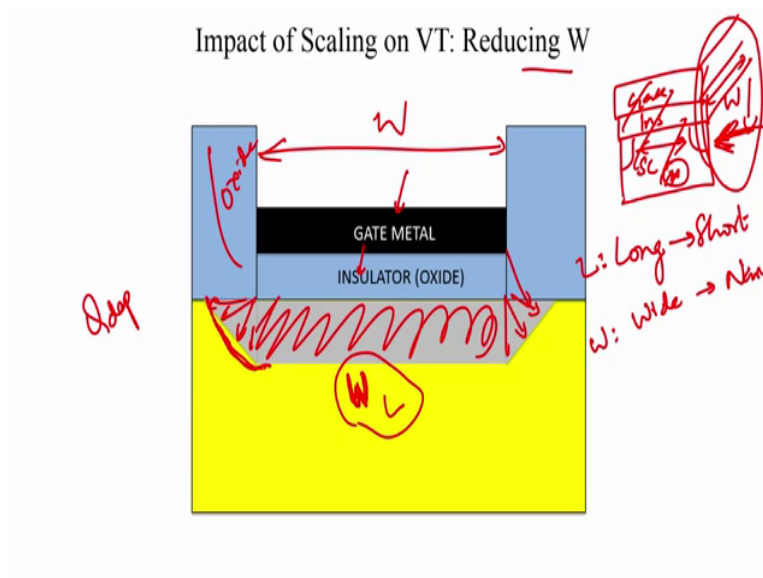


So now, what about the impact of scaling on the of the channel width ok; so, we have your we have the MOSFET which is got the it is the gate insulator semiconductor and you now have, we have discussed the impact of channel length, but what about reducing channel width. So, for reducing channel length, we say it is a long channel device to a short channel device, but for channel width, we say it is a wide device to a narrow device.

Now, let us take the cross section of a wide device ok. We are going to look at this device from this side ok. We have always been looking at the device inside the plane. We have always looked at this cross section, but now let us start looking at that cross section ok. So, that cross section is shown here ok. So, here I have taken a cross section right, in the centre of the device. So, you cannot see the source and drain. We have a bypassed area somewhere, in the middle of the channel and you have your gate metal you have the insulator, that is a semiconductor and that is the depletion region and you have oxides on all these side ok. On the side walls it is all insulator and the source is somewhere, outside the screen and the drain is somewhere inside the screen.

So, we are looking the channel the channel lies in and out of this plane. So, we are looking through the device in this direction and therefore, this is the channel width of the device or this is the intended channel width of the device. Now, what is the impact of reducing this channel width? As we scale down, impact of reducing the channel width on the threshold voltage of the MOSFET, so, if the channel width; so, let us look at the depletion profile.

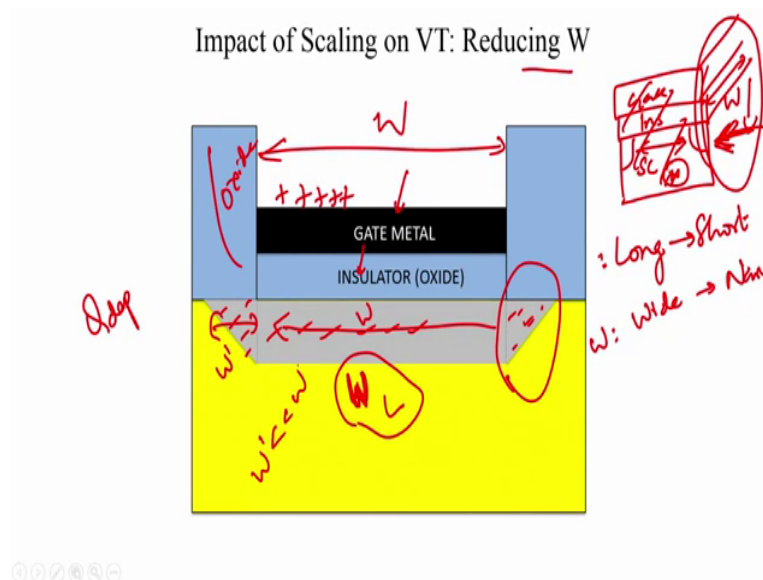
(Refer Slide Time: 31:55)



So, once again the answer has got to do with the depletion charge. So, total depletion charge. So, if you look at this depletion profile, yes, most of this depletion, if you look at the total depletion charge all depletion charge spans over a length of W and of course, the channel length to say, L and therefore, W into L is what is depleted, but then, because of the fringing in the field ok.

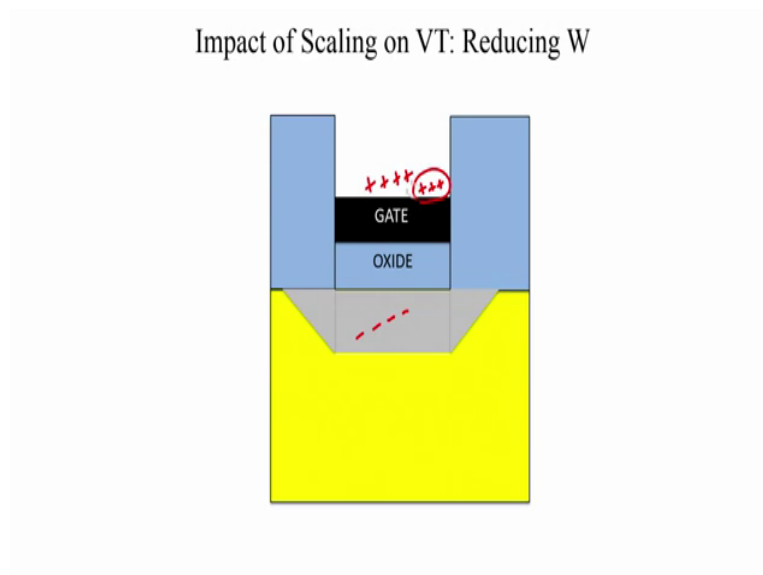
So, you do have some field lines ok, heading towards the sides and therefore, due to this fringing of the field you are going to end up, whether you like it or not we are going to end up depleting some excess charge on the sides and this excess charge is not really contributing too much to the behavior of the MOSFET. But, it is simply that if you add if you add some charge in the gate you are going to end up not only depleting this region, but you are also going to unintentionally deplete the regions on the sides, because of the fringing of the fields.

(Refer Slide Time: 32:39)



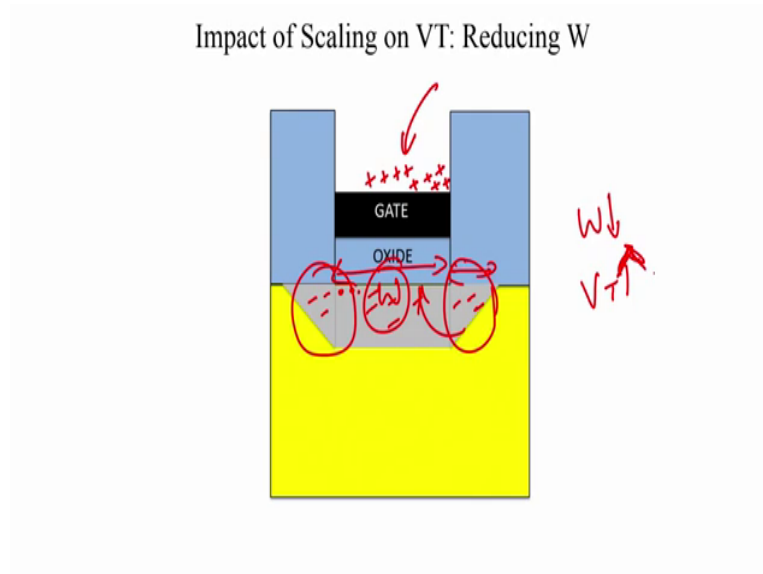
Now, for a very wide device this extra work done, really does not make any difference ok. It is a very small percentage of this total channel width. So, this W_{dash} is a very small percentage of W , but for a short channel device it becomes a very significant portion ok.

(Refer Slide Time: 33:19)



So, let us say let us say if, let us say, we intended to add some positive charge on the gate and quickly deplete this region and after that, when we add some more extra charge, we expect inversion ok, or the definition of V_t , but now since the extra charge so now, since, whatever charge is being added there is also compensated for, by these regions on the sides right.

(Refer Slide Time: 33:49)



So, I still need to add more charge and once again, that is also compensated for this and only when you add a lot more charge will the depletion region finally have been fully depleted and then, allow inversion. So, you need to add a lot more charge in the gate, because these regions on the sides are helping or, helping the semi conductor respond to this applied charge. And therefore, you need to do a lot more work, before you get inversion and therefore, reducing the channel width increases the V_t and it is again a percentage question right.

So, this is my channel width W , these regions have now become significant compared to this W and therefore, the V_t actually increases. And, if you want to see the mathematics of it, you can say that, let us say, what is the difference between a wide channel device and a narrow channel device the wide channel device is got what is a total depletion charge. Again, we are going to measure the charge in coulombs and not per unit area or volume.

(Refer Slide Time: 34:41)

Impact of Scaling on VT: Reducing W

$$V_{T-Wide} - V_{T-Narrow} = \frac{-Q_{dep-Wide}^{total} + Q_{dep-Narrow}^{total}}{WLC_{ox}} = \frac{qN_A}{C_{ox}} \left[\frac{WLx_d}{WL} - \frac{WLx_d + \pi x_d^2 L / 2}{WL} \right]$$

$$= -\frac{qN_A x_d}{2C_{ox}} \left[\frac{\pi x_d}{L} \right]$$

Reducing Channel Width => VT increases

(Refer Slide Time: 34:57)

Impact of Scaling on VT: Reducing W

$$V_{T-Wide} - V_{T-Narrow} = \frac{-Q_{dep-Wide}^{total} + Q_{dep-Narrow}^{total}}{WLC_{ox}} = \frac{qN_A}{C_{ox}} \left[\frac{WLx_d}{WL} - \frac{WLx_d + \pi x_d^2 L / 2}{WL} \right]$$

$$= -\frac{qN_A x_d}{2C_{ox}} \left[\frac{\pi x_d}{L} \right]$$

Reducing Channel Width => VT increases

The total depletion charge in a wide channel device the total depletion charge is your $q N_A x_d$ into W into L and that divided by the total gate capacitance is nothing, but C_{ox} into $W L$, what is the total depletion charge in a narrow channel device. Now, the side regions cannot be ignored. So, even here the side regions are present, but W plus anything on the side is approximately W . This is this excess, charge is very-very miniscule compared to the total charge in the gate, but now it is not ok.

(Refer Slide Time: 35:43)

Impact of Scaling on VT: Reducing W

$$V_{T-Wide} - V_{T-Narrow} = \frac{-Q_{dep-Wide} + Q_{dep-Narrow}}{WLC_{ox}} = \frac{qN_A}{C_{ox}} \left[\frac{WLx_d}{WL} - \frac{WLx_d + \pi x_d^2 L / 2}{WL} \right]$$

$\frac{Q_{dep-wide}}{C_{ox}} = \frac{qN_A x_d WL}{C_{ox} WL}$

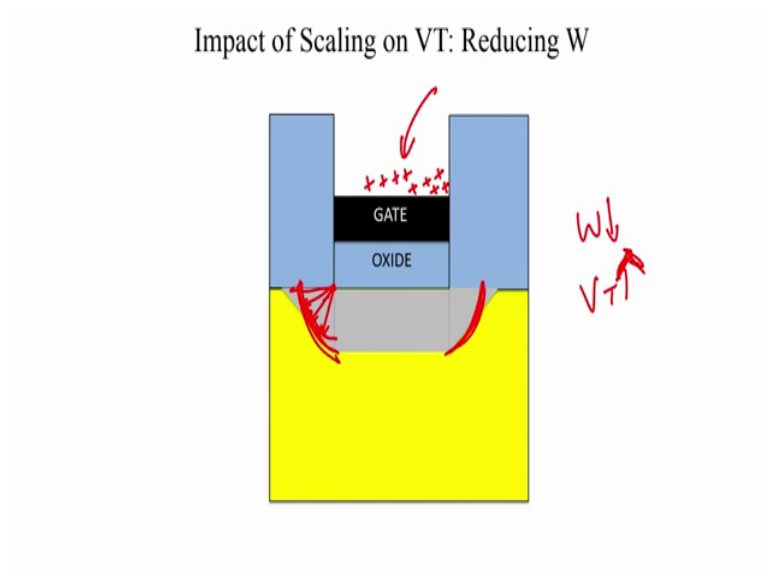
 $\frac{Q_{dep-narrow}}{C_{ox}} = \frac{qN_A x_d (WLx_d + \frac{\pi x_d^2 L}{2})}{C_{ox} WL}$

Reducing Channel Width \Rightarrow VT increases

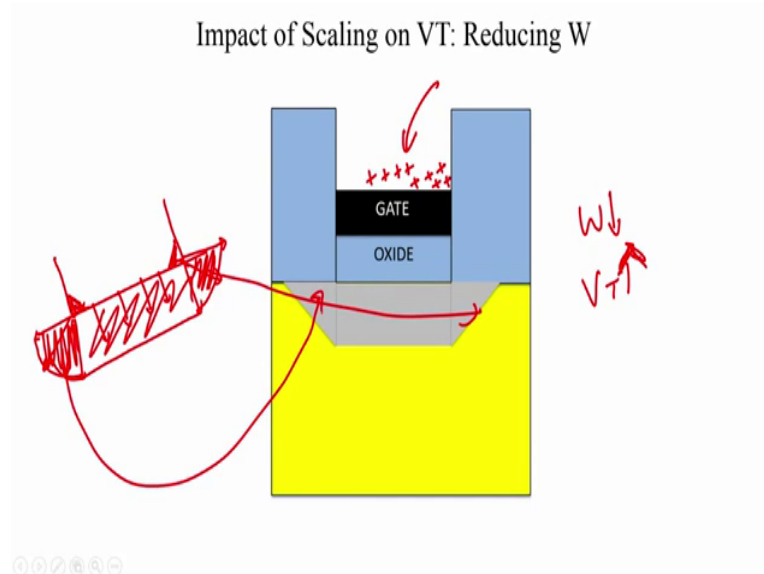
$V_{T-wide} < V_{T-nar}$
 $L \downarrow \quad V_T \downarrow$
 $W \downarrow \quad V_T \uparrow$

So, in this case, it is going to be $q N_A x_d$ by $C_{ox} WL$ into $W L x d$, which is this region which is your regular which is the regular region that you have already, that you need to deplete. But we are now depleting some extra regions, and we are going to model it model this extra region not like this little triangles, but you know model, it along these fringing lines the fringing electric field ok. So, that is it is going to have a more of a circular profile. So, you are going to have the circular arcs. So, although I have shown it as, although I have shown it, as you know, straight edges just for this just, because it is easier to draw.

(Refer Slide Time: 36:13)



(Refer Slide Time: 36:29)



You have a gate you have an insulator here and this insulators are sharp if you look at this edge of this insulator. It is a sharp a corner and this corner is going to have a very strong fringing field. So, that is the way the field lines are going to be present it is going to have a very strong field. And therefore, this is the depletion region, it is you can say it is more like a sector of a circle and then here you have this region that is depleted.

So, that corresponds to this and this, corresponds to this. So, if you take that little accuracy if bringing, that little accuracy in the picture you will find that, we use the area of that little sector to define this excess area being depleted that divided by $W L$, that divided by C_{ox} into $W L$ is the is this term. So, what is the difference between these two terms? The difference will now be a negative number, which implies that my threshold voltage for the wide channel device is going to be less than the threshold voltage for a narrow channel device.

So, in summary as we scaled down L the V_T also decreased is the transistors, turning on faster and why is that, because the PN-junction is helping is aiding the depletion. But, if you scale down W the V_T increases and why is that it is, because you have to do some extra work in depleting some regions on the sides of the channel ok. So, I hope this qualitative and quantitative explanation helps out, but if there are any questions do feel free to send me an email.