



Laboratoire IMS Bordeaux

université
de BORDEAUX

BORDEAUX
INP



  | IMS Bordeaux



Emerging technology and new computing paradigms: How to reinvent the value chain?

Cristell Maneux¹, Marina Deng¹, Chhandak Mukherjee¹, David Atienza², Jens Trommer³, Oskar Baumgartner⁴, Guilhem Larrieu⁵ and Ian O'Connor⁶

¹University of Bordeaux,FR ; ²École Polytechnique Fédérale de Lausanne (EPFL), CH; ³Namlab gGmbH, DE; ⁴Global TCAD Solutions, AT; ⁵LAAS – CNRS, FR; ⁶Lyon Institute of Nanotechnology, FR



Outline

- Contexts
- Characterization challenges
- Modelling challenges
- Circuit design challenges

Contexts

○ 4th industrial revolution

- Unprecedented growing demand for neural networks (NNs)
- **But**, technology solutions still relying on transistors inherited from Moore's Law, optimized for von Neumann machines.

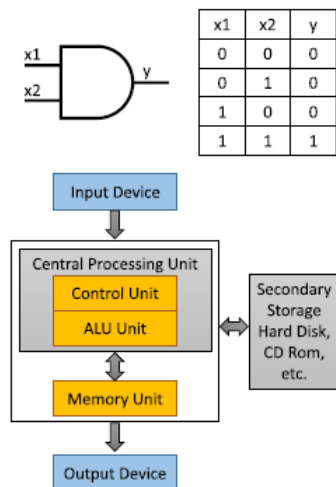
○ Energy efficiency of von Neumann based processors

- limited by data transfer between memory and computing cores implemented in 2D integration schemes.
- **Crippling limitation** for NNs efficiency.

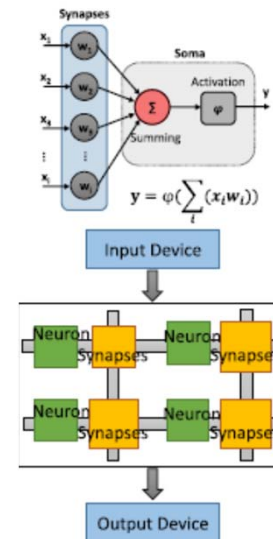
Contexts

Machine's concept : von Neumann based

Von Neumann architecture

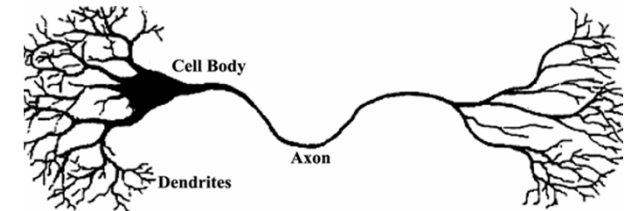
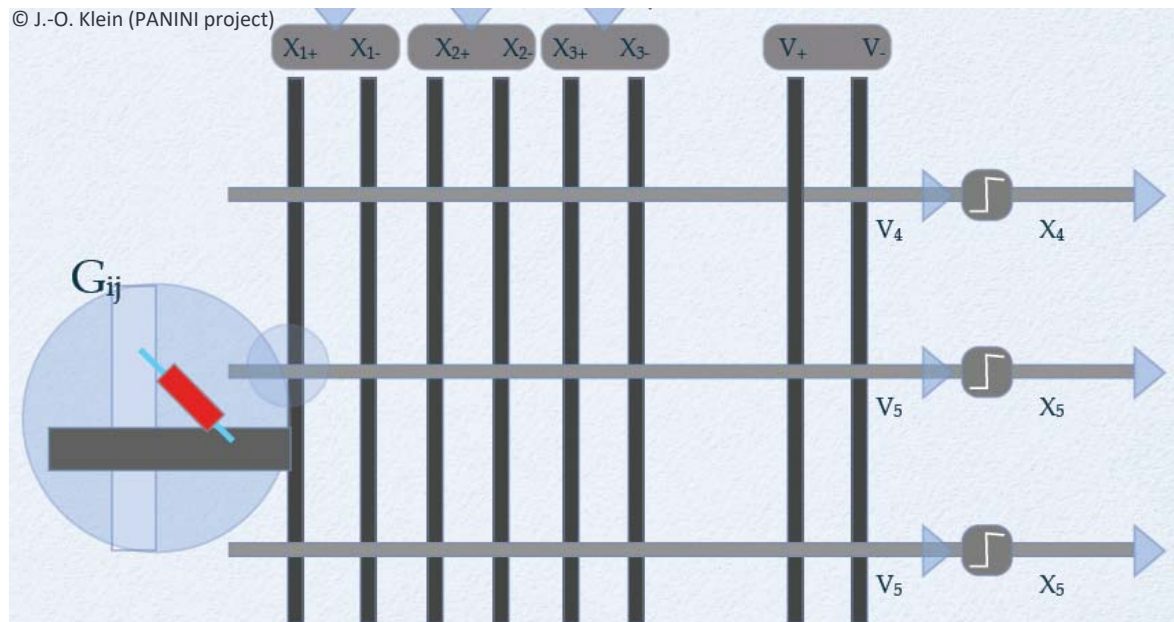


Neuromorphic architecture



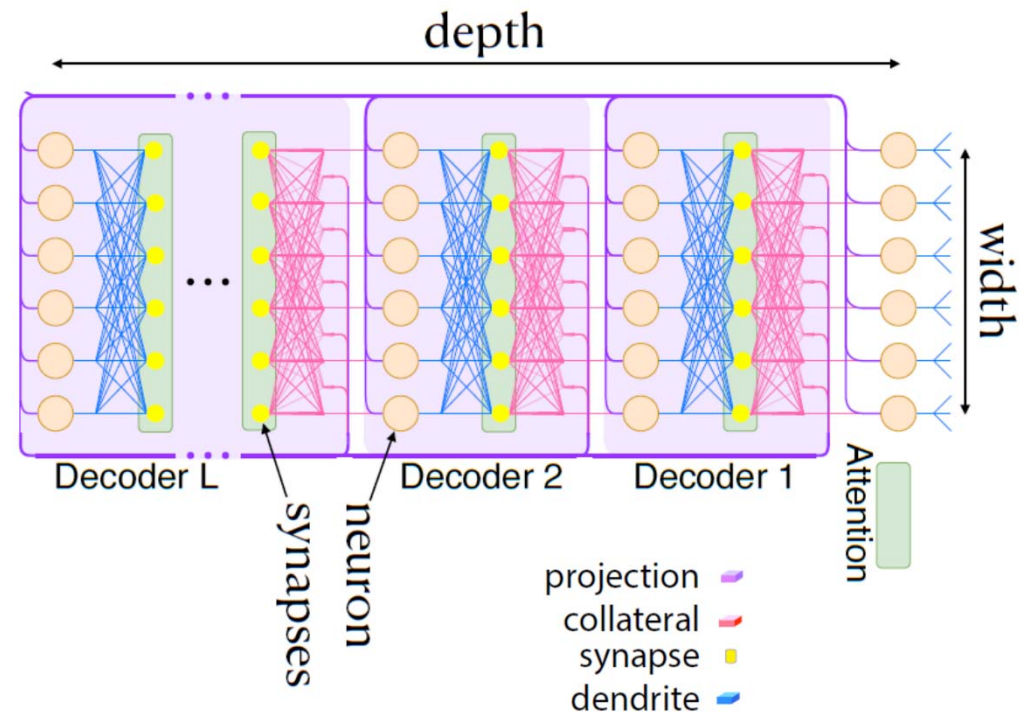
Contexts

Machine's concept: Neural Networks



Contexts: Architecture of Neural Networks

- Dense **local connections**
 - feedforward or recurrent
- Sparse **global projections**
 - residual or feedback
- Size = depth × width
 - Depth = Number of stacked layers
 - Width = Number of neurons in a layer



Contexts: 2D Neural Networks

○ "Compute-in-memory"

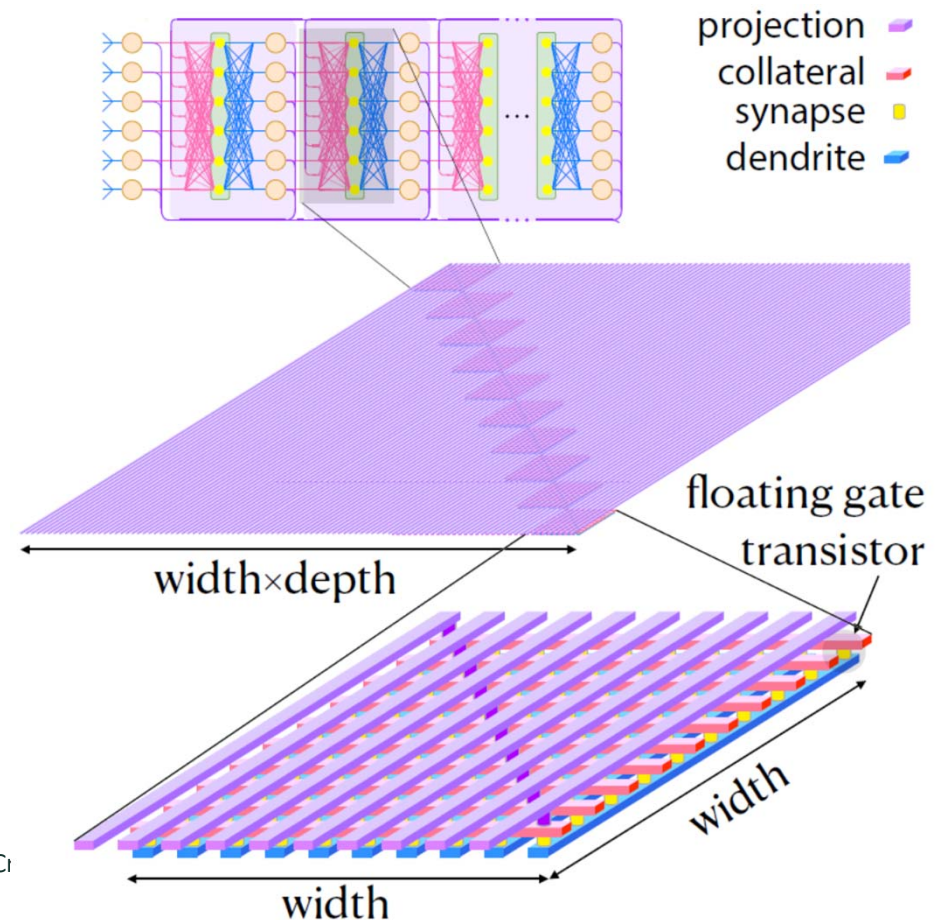
○ $\text{work} \propto \text{distance} = \text{width} \times \text{depth}$

○ $\text{signals} \propto \text{size} = \text{width} \times \text{depth}$

○ $\text{energy} = \text{work} \times \text{signals} \propto \text{size}^2$

○ $\text{area} \propto (\text{width} \times \text{depth})^2 = \text{size}^2$

○ thermally *viable*



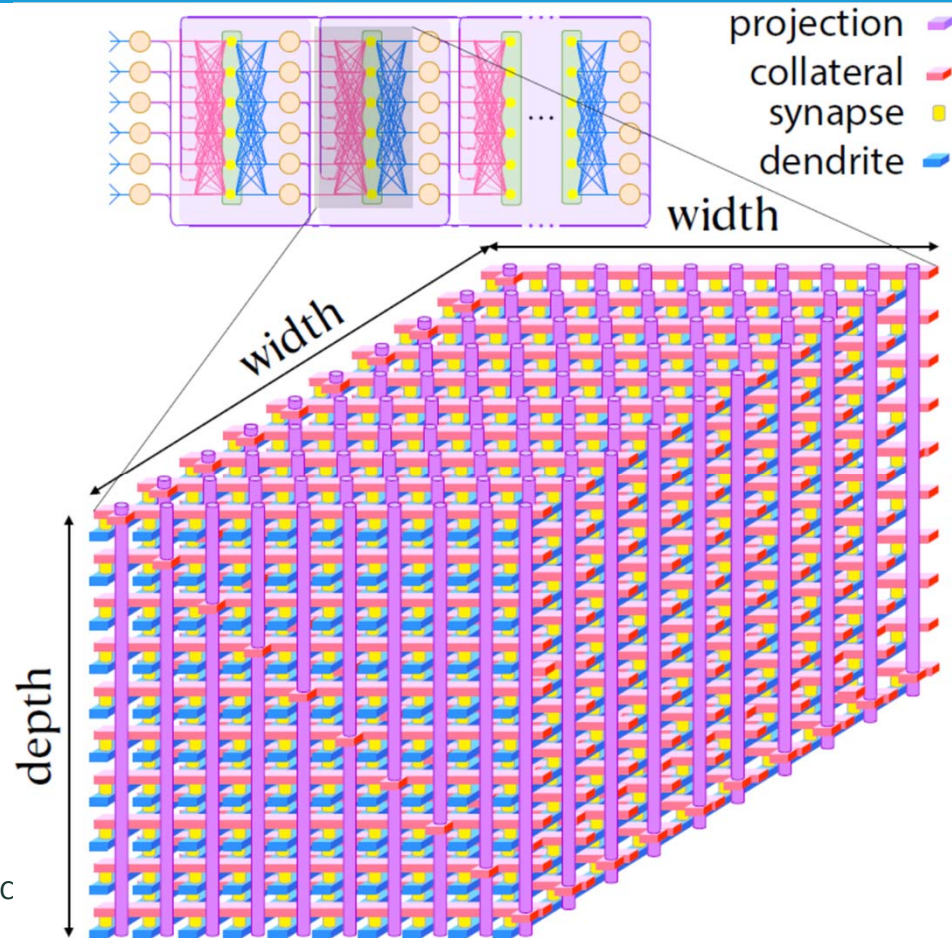
Contexts: 3D Neural Networks

- For width \gg depth

- work \propto distance \approx width
- signals \propto size = width \times depth
- energy = work \times signals
= width² \times depth

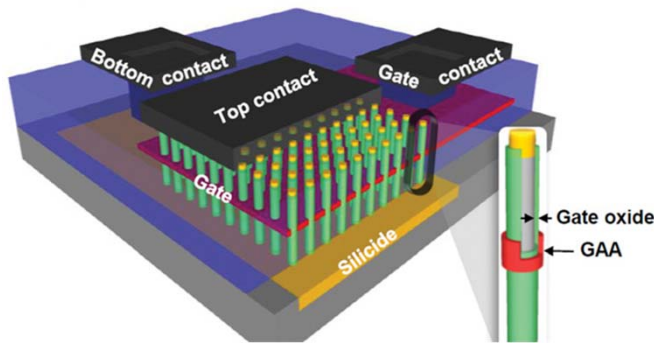
- And depth \propto width

- energy = size^{1.5}
- area \approx width² \propto size
- Thermally viable?

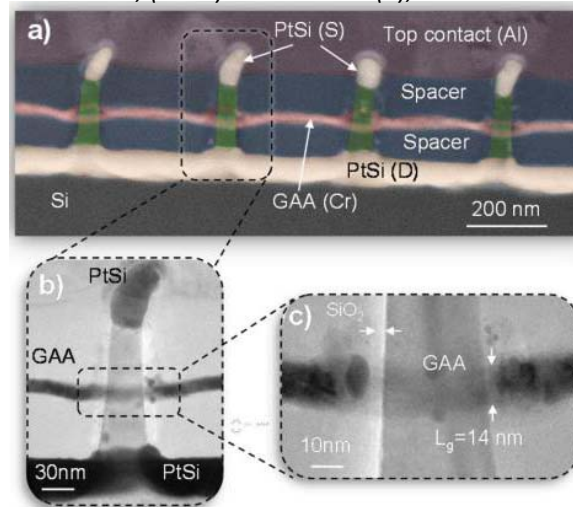


Context: Conclusion

- To save energy, we need :
 - A 3D implementation ...
 - Dedicated devices for 3D stacked NNs
- How could we do this?



G. Larrieu, (2013) Nanoscale 5 (6), 2437-2441.



Key features of the technology

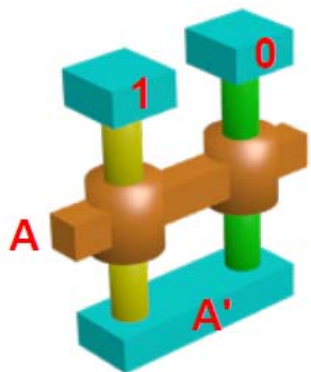
- Vertical Gate-all-around device
- Scaled gate length ~ 14 nm
- Symmetrical S/D silicided contacts that ensure **junctionless operation**

Outline

- Contexts
- Characterization challenges
- Modelling challenges
- Circuit design challenges

Vertical Nanowire FET: Characterization challenges

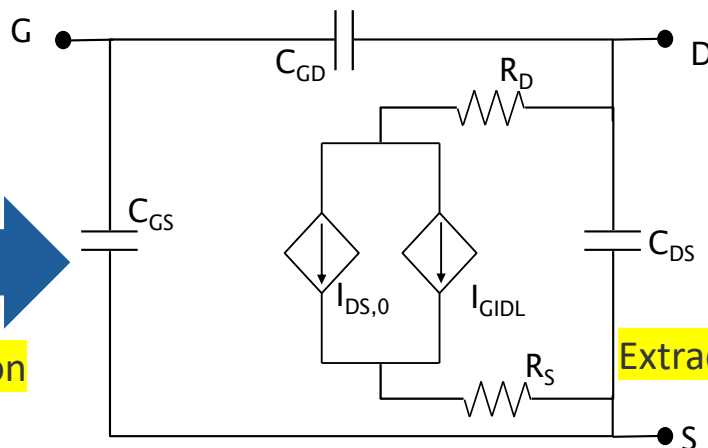
○ Measurement accuracy



INV1-SL-NP

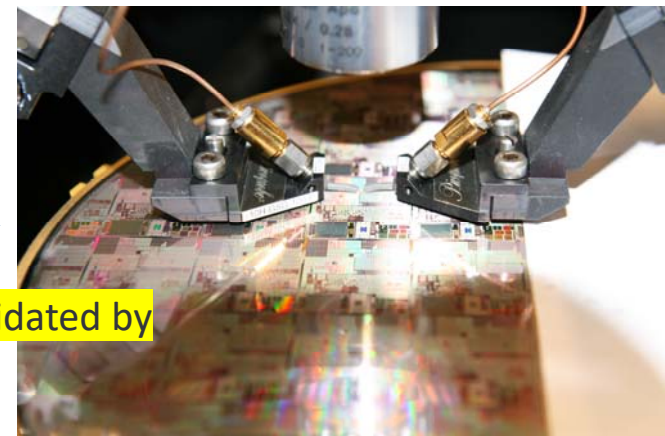
Circuit/system design

Relies on



Accurate compact models

Extracted from/validated by



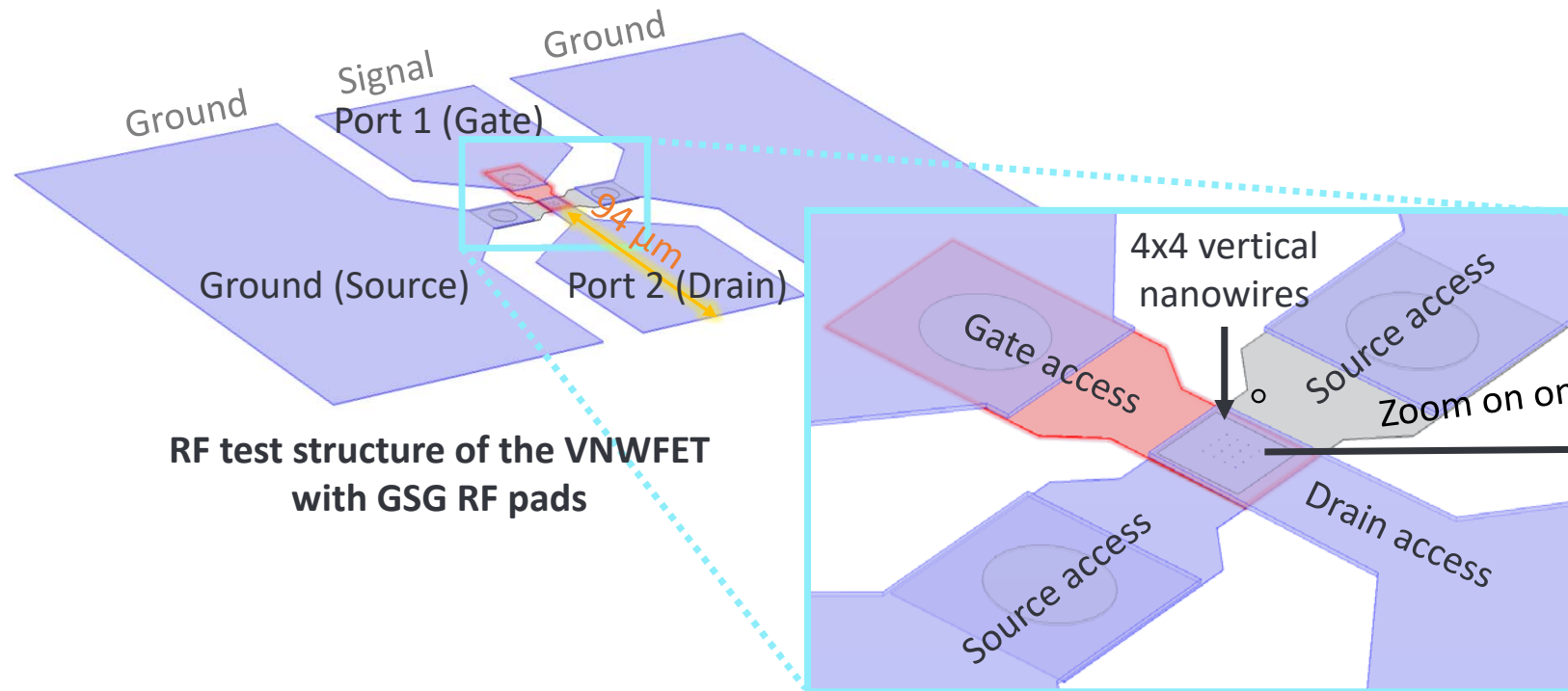
On-wafer device measurements
DC + S-parameters



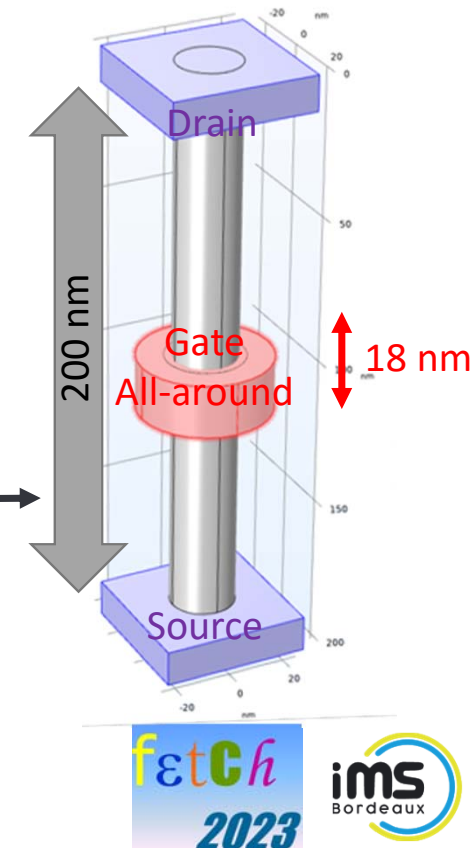
Accurate S-parameter measurements for nanoscaled devices=big challenge !!!

Vertical Nanowire FET: Characterization challenges

○ From test structure to intrinsic device



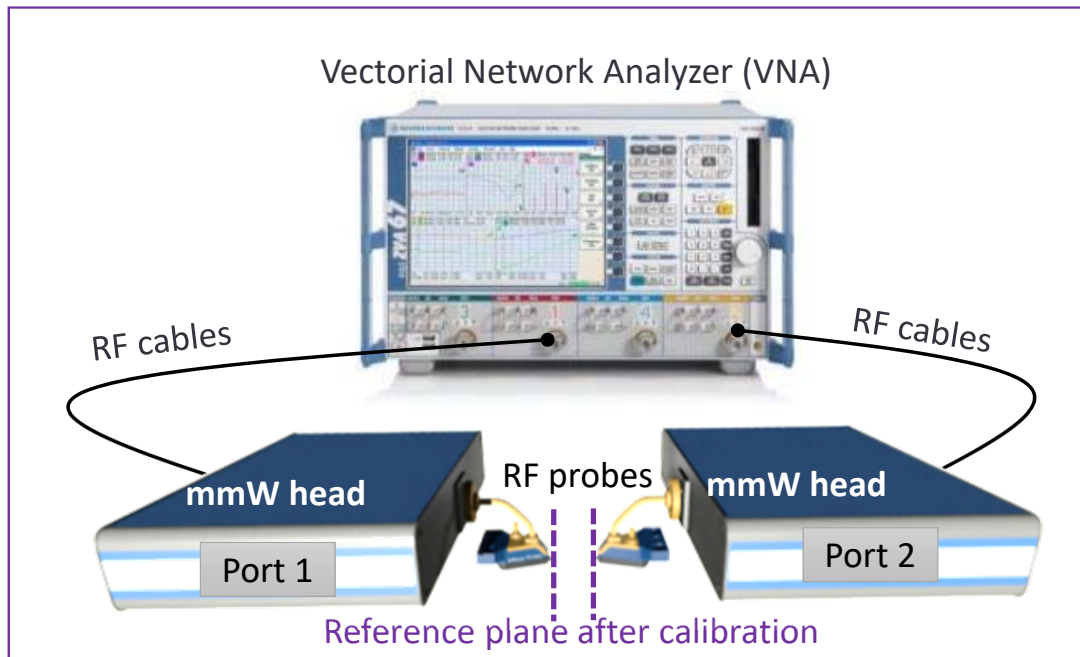
Single nanowire constituting the VNWFET



Vertical Nanowire FET: Characterization challenges

Calibration & De-embedding Issues

- Standard procedure=2-step calibration



1

Off-wafer
calibration on ISS

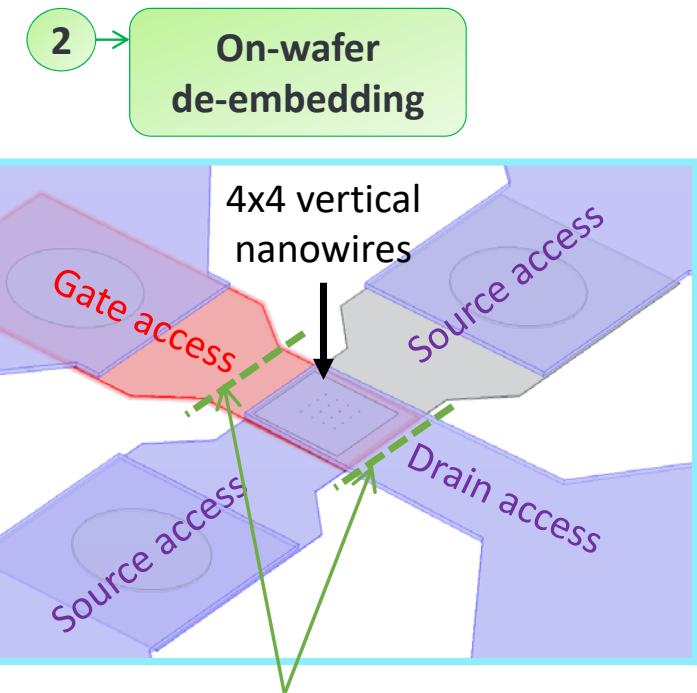
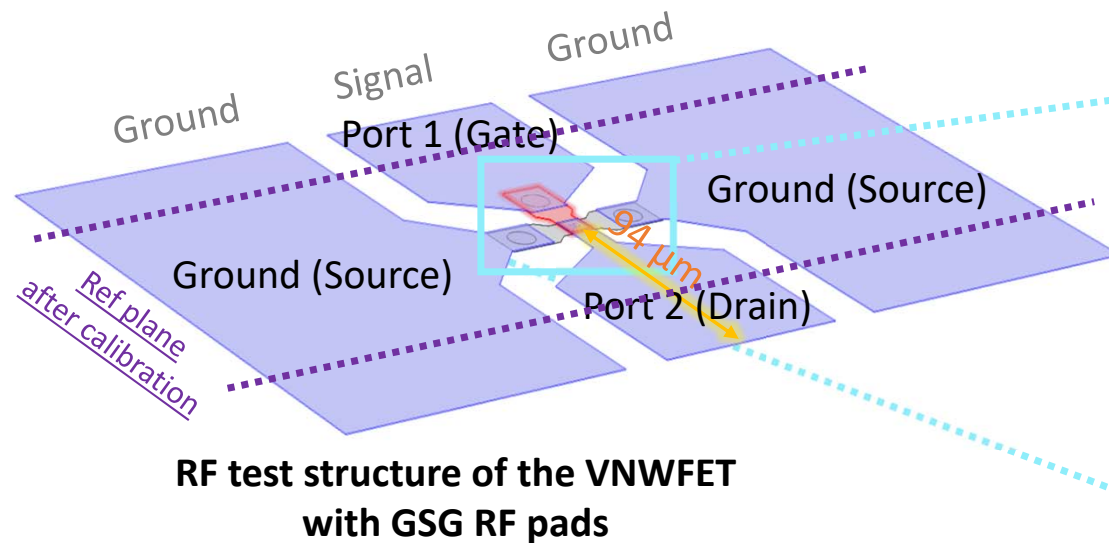


Impedance Standard Substrate
Source=Cascade Microtech

Vertical Nanowire FET: Characterization challenges

Calibration & De-embedding Issues

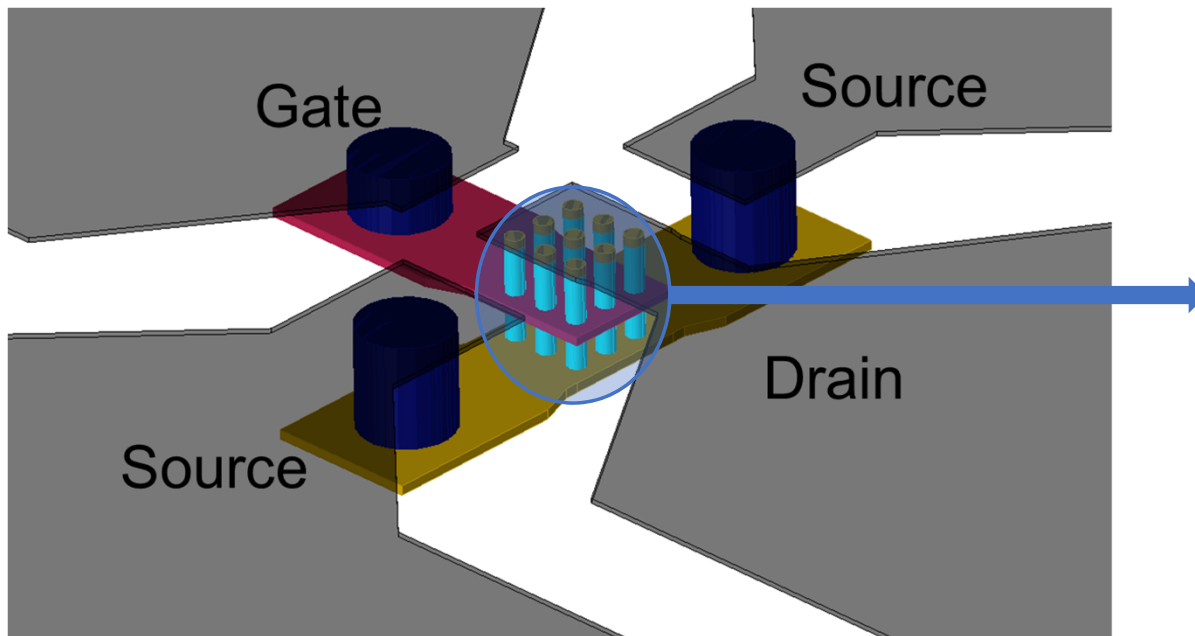
- Standard procedure=2-step calibration



VNWFET = 500x smaller than the accesses

→ De-embedding step is critical

Which de-embedding method?



Inside: the intrinsic device,
P type 81 **vertical** nanowires
with 17 nm of diameter each

Outside: **3D Interconnects** and
pads necessary for electrical
contact (DC and RF pads)

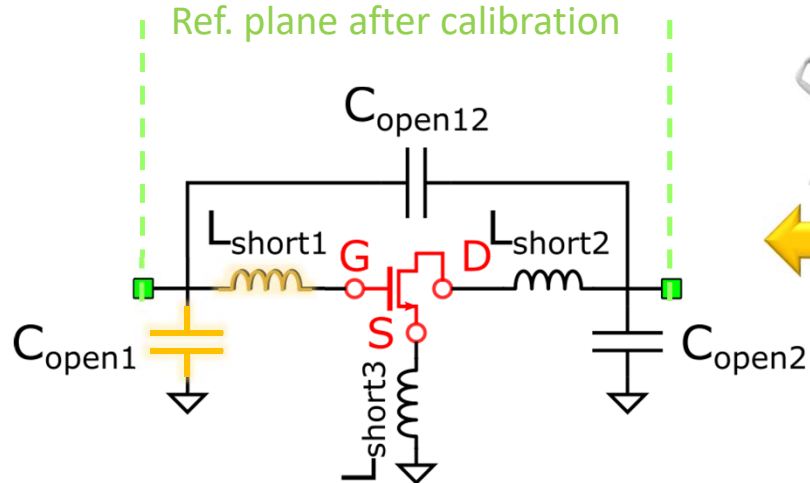
Which de-embedding method?

○ De-embedding removes errors due to device accesses

↳ Based on an equivalent circuit of the transistor accesses

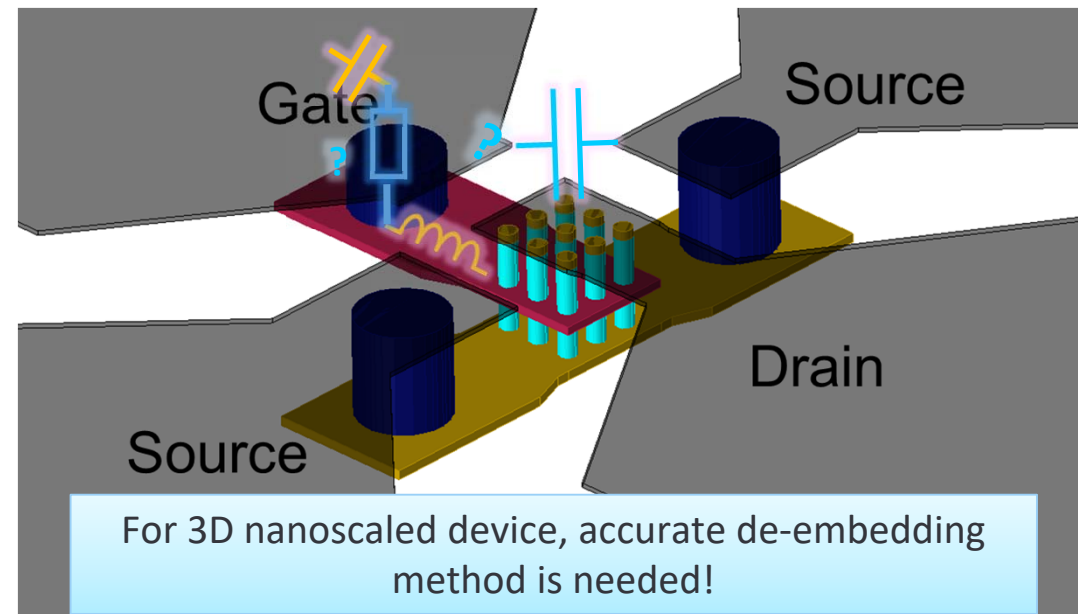
Small Signal Equivalent Circuit: SSEC

Ref. plane after calibration



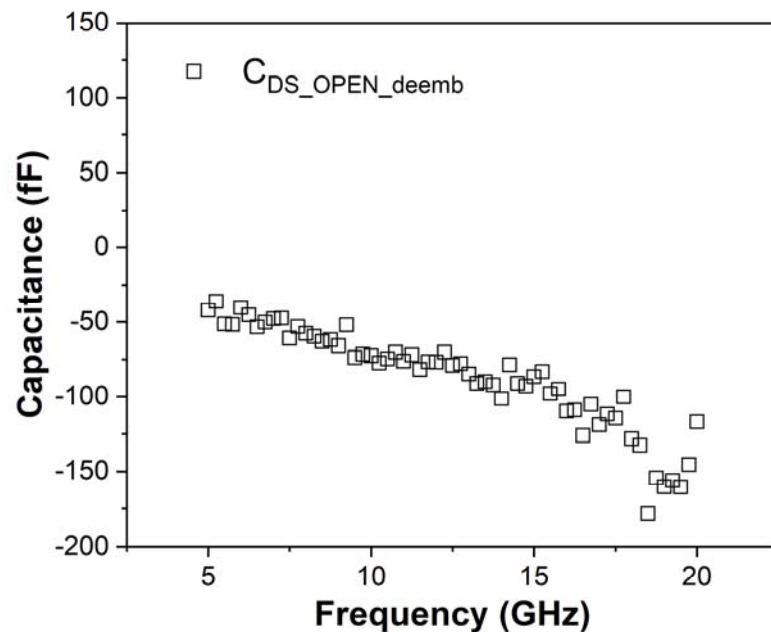
OPEN-SHORT DE-EMBEDDING

(conventionnal method)



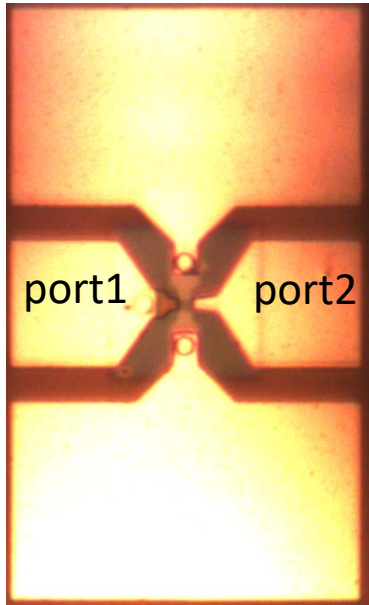
Which de-embedding method?

Open de-embedding



- Unphysical capacitance value
- New de-embedding method necessary

De-embedding method for VNWFETs



OPEN

Measure-
ment

- S-parameter measurement up to 40 GHz
- DUT = Open test structure and transistor

EM
simulation

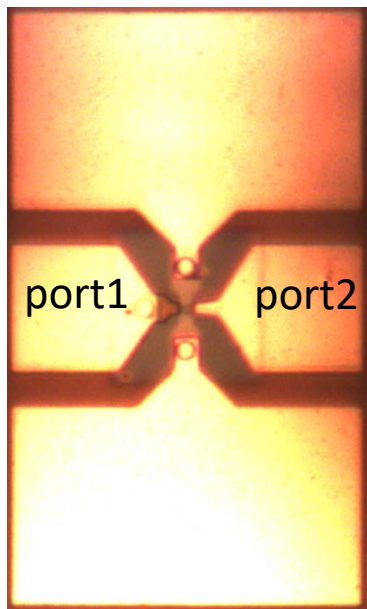
- Simulator = ADS-Momentum
- Definition of the electrical properties for each layer used in the VNWFET process

Electrical
simulation

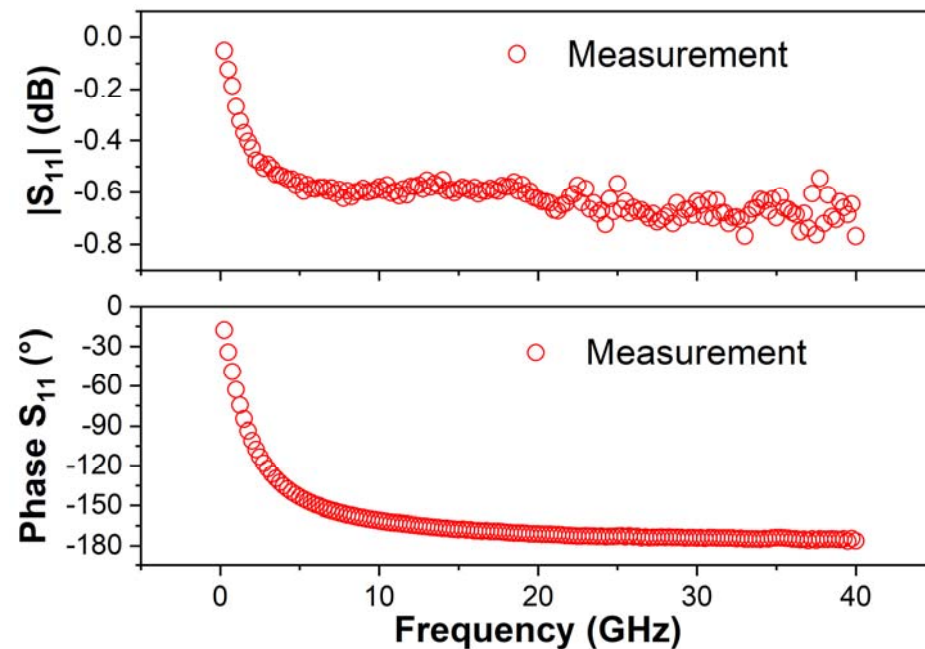
- SSEC of the parasitic network
- De-embedding of VNWFET

S-parameters measurements

- S-parameters on-wafer measurements up to 40 GHz of the OPEN structure

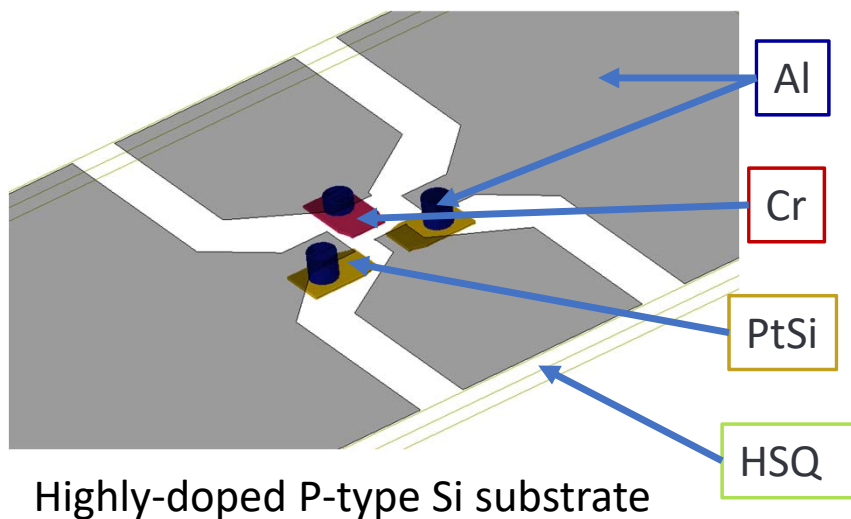


OPEN

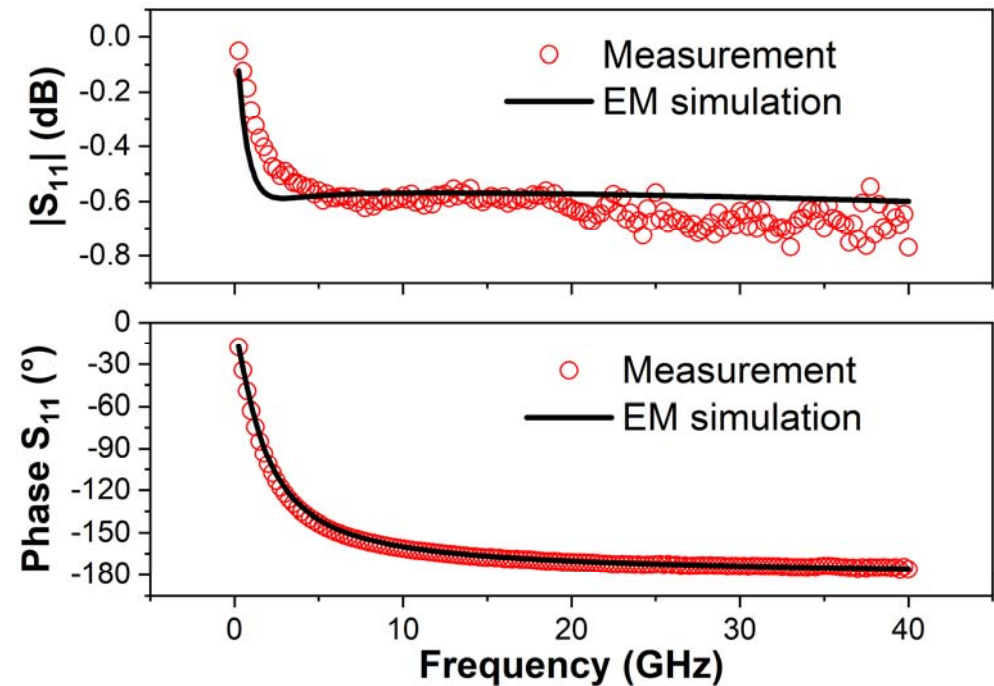


EM simulation

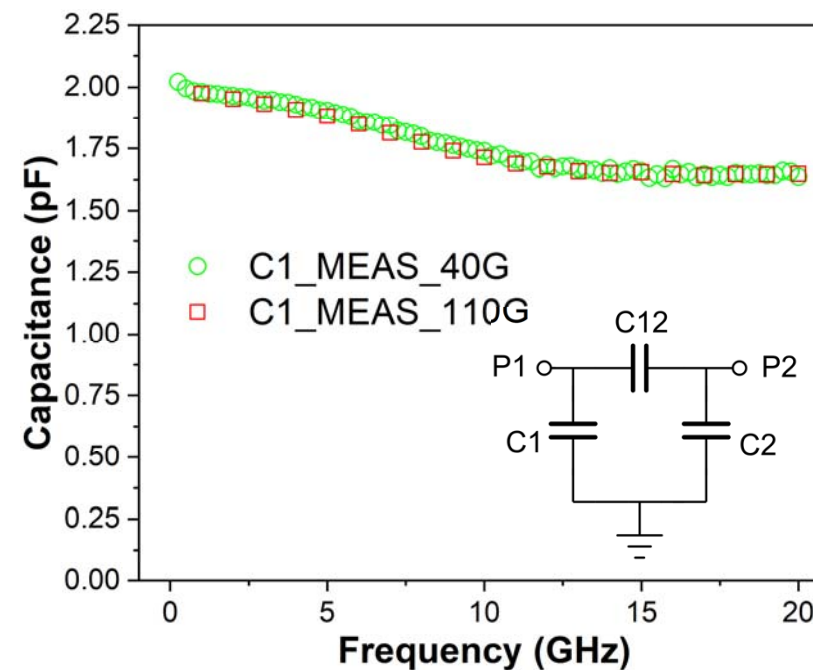
○ S-parameters from EM simulation up to 40 GHz of the OPEN structure



OPEN ADS 3D VIEW



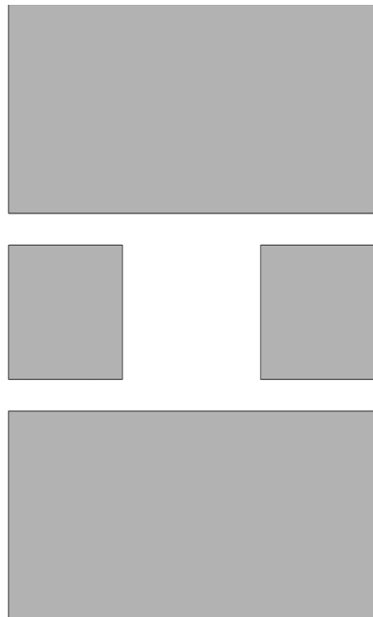
Open Capacitance extraction



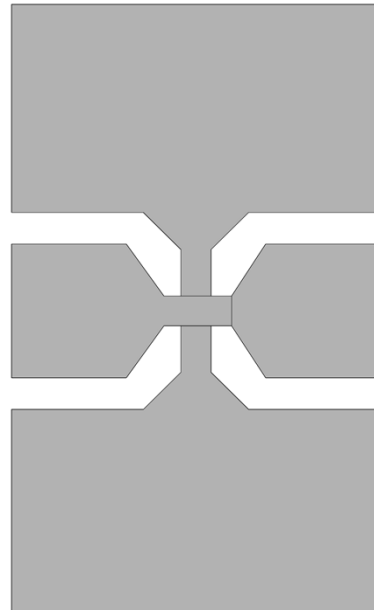
○ Open structure as a π -network? Investigate a more complex circuit!

SSEC of the parasitic network

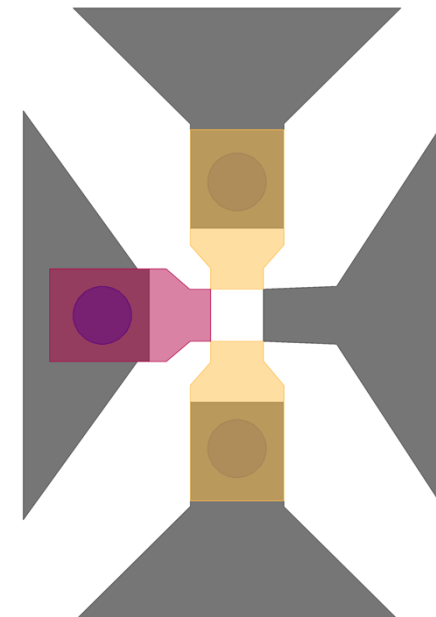
○ EM simulation for virtual structures



OPEN PAD



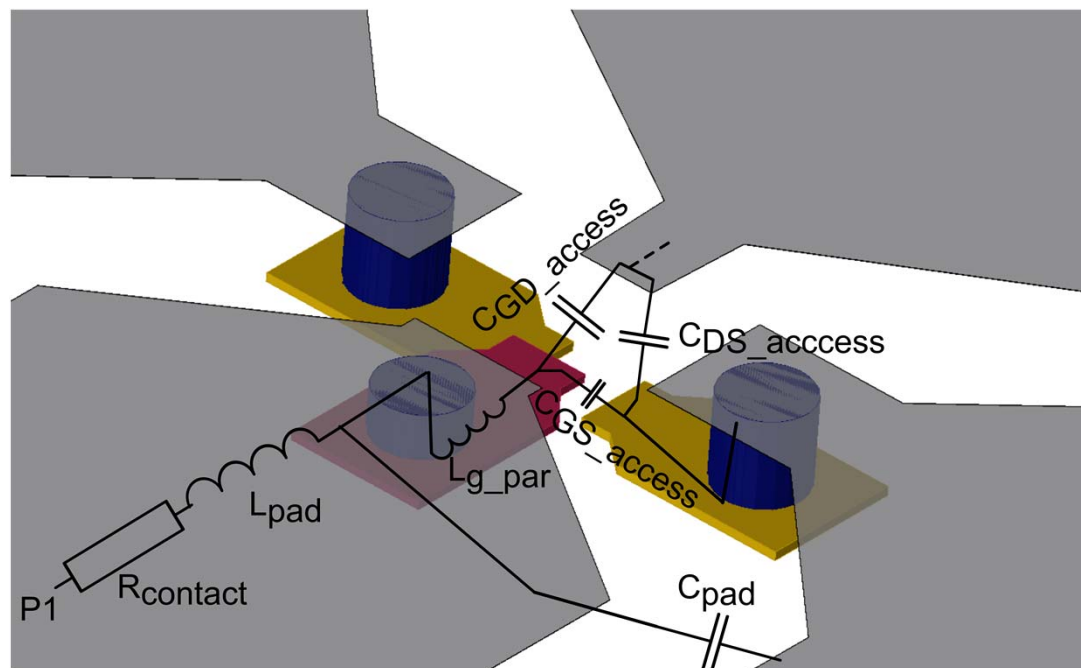
SHORT PAD



OPEN DEVICE

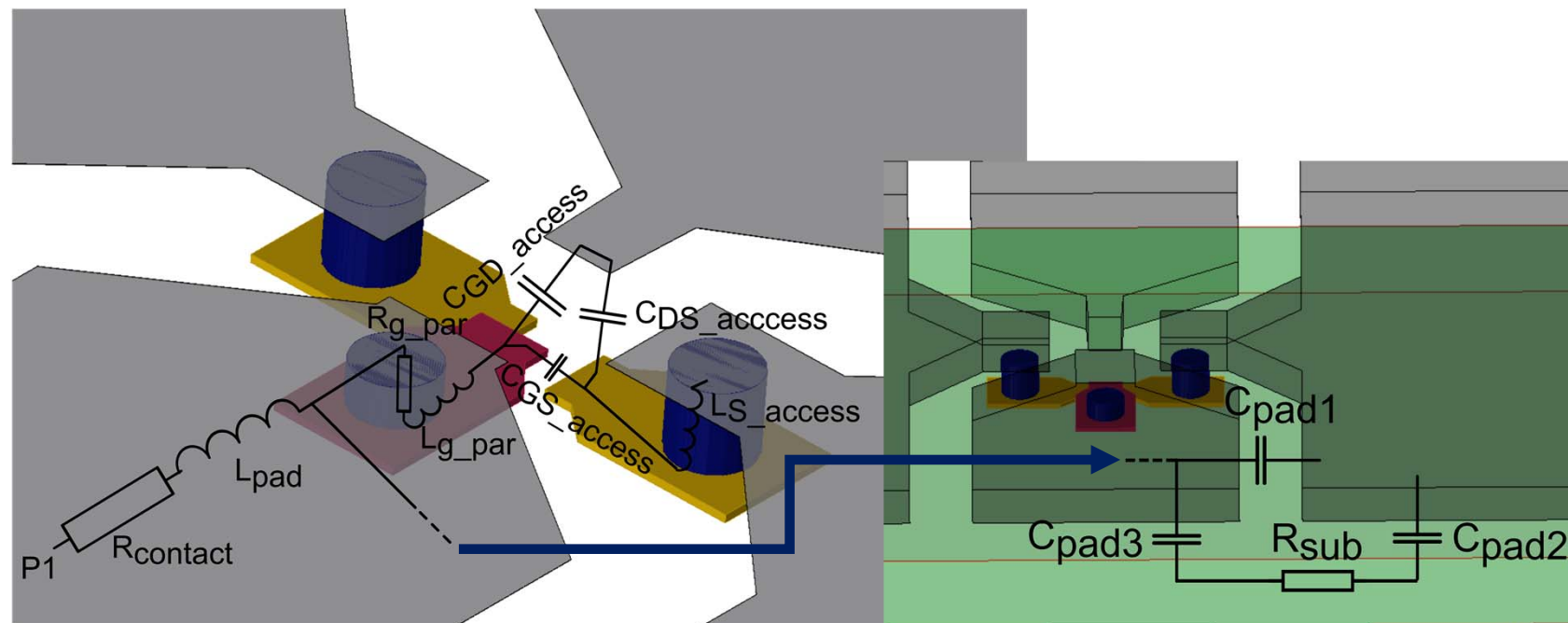
SSEC of the parasitic network

○ Distribution of the lumped elements: MODEL1



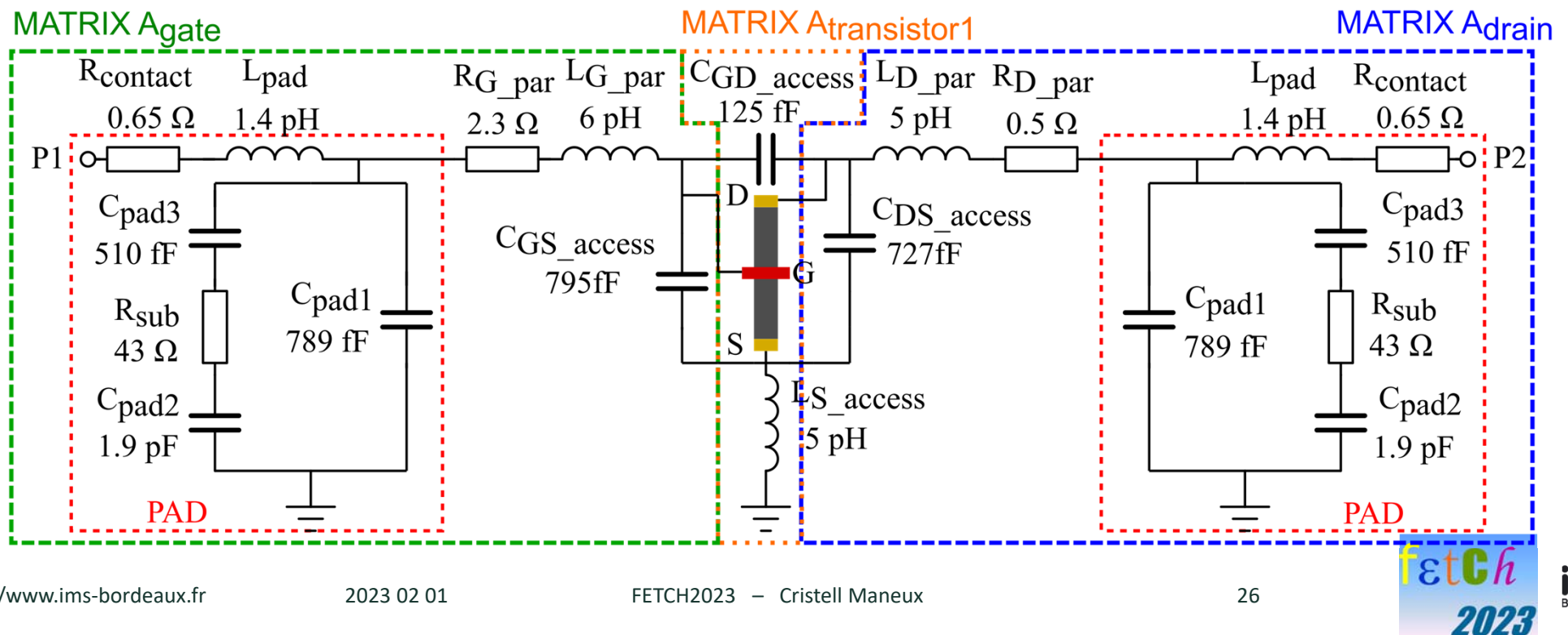
SSEC of the parasitic network

○ Distribution of the lumped elements: MODEL2



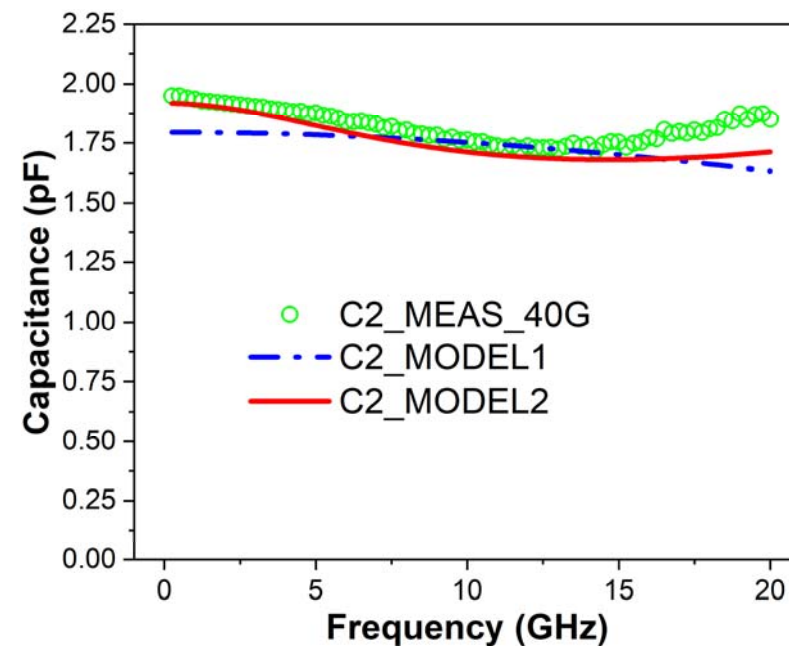
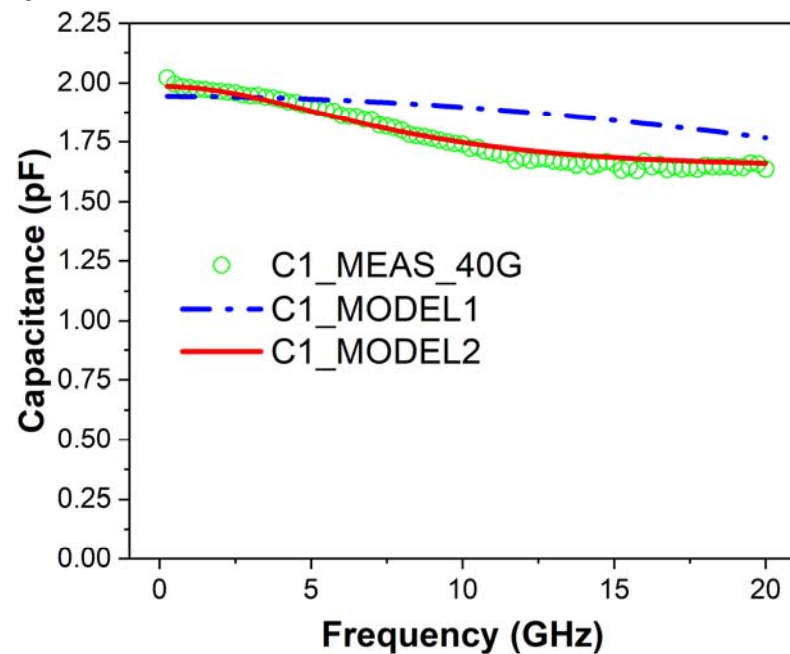
SSEC of the parasitic network

- Final small signal equivalent circuit of the parasitic network: MODEL2



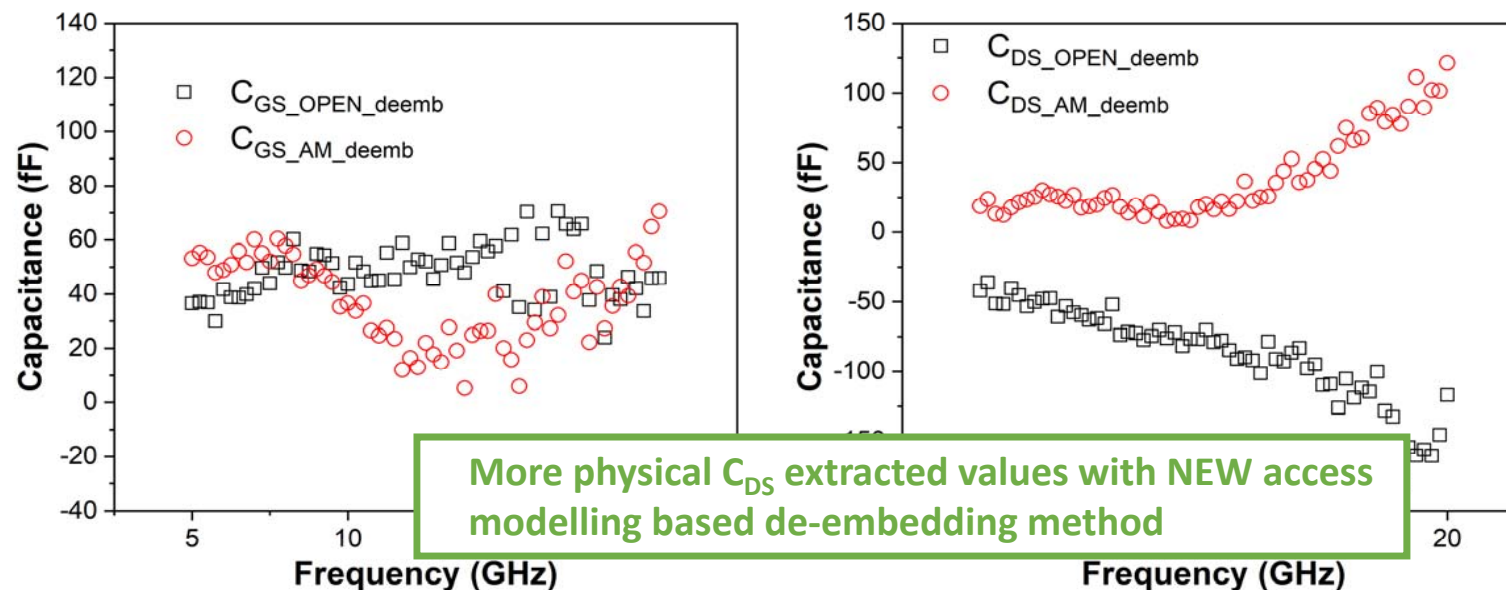
SSEC of parasitic network

- Capacitance from the small signal equivalent circuit of the open structure up to 20 GHz



De-embedding results

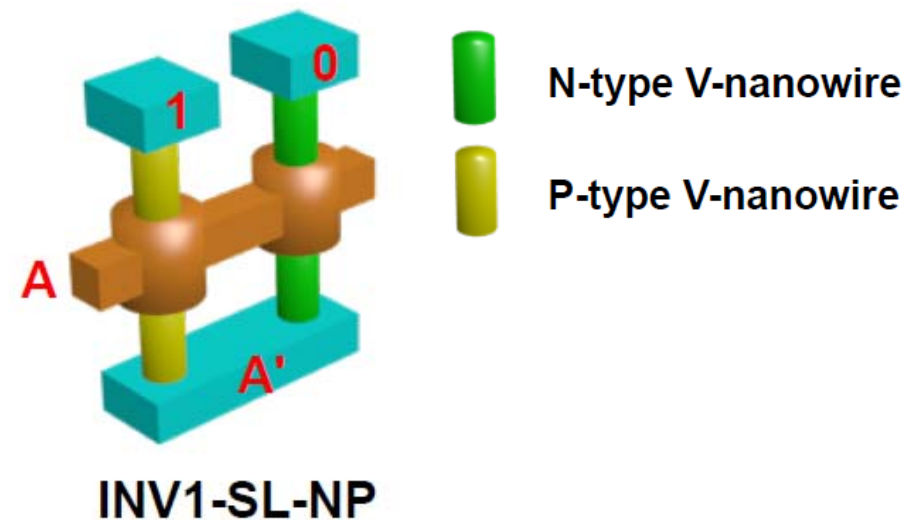
- Comparison between open and access modelling de-embedding for C_{GS} and C_{DS} extraction



VNWFET Characterization

○ Reliable measurements for VNWFET model parameter extraction is OK!

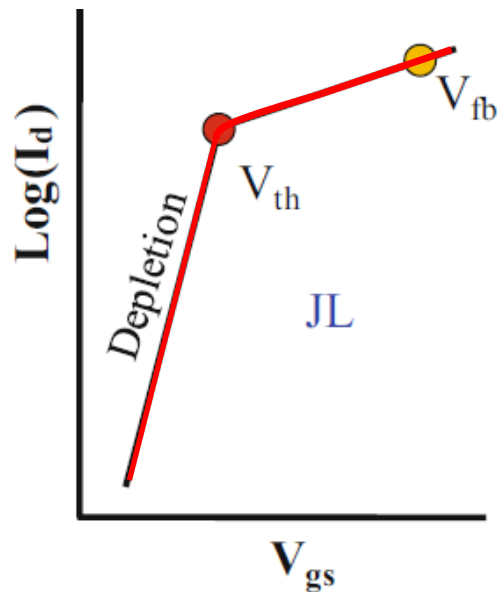
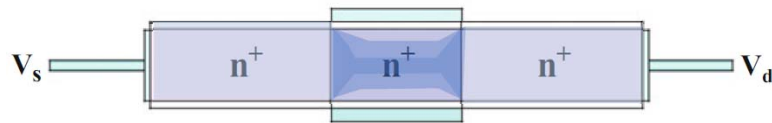
➔ Ready for compact modelling!



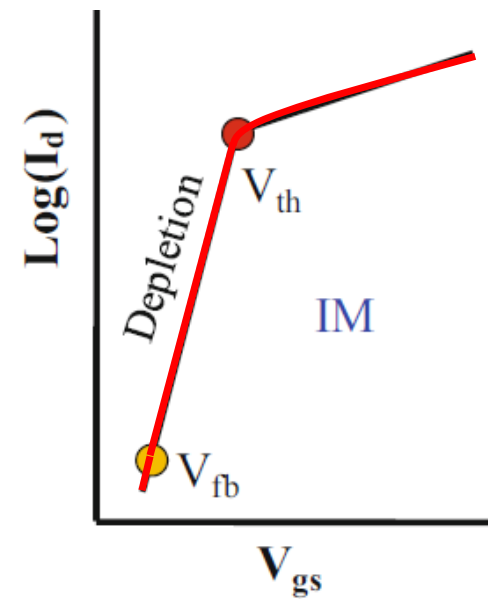
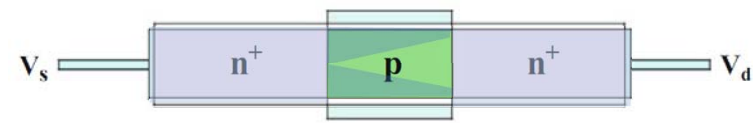
Outline

- Contexts
- Characterization challenges
- Modelling challenges
- Circuit design challenges

Junctionless vs. Classical MOSFETs

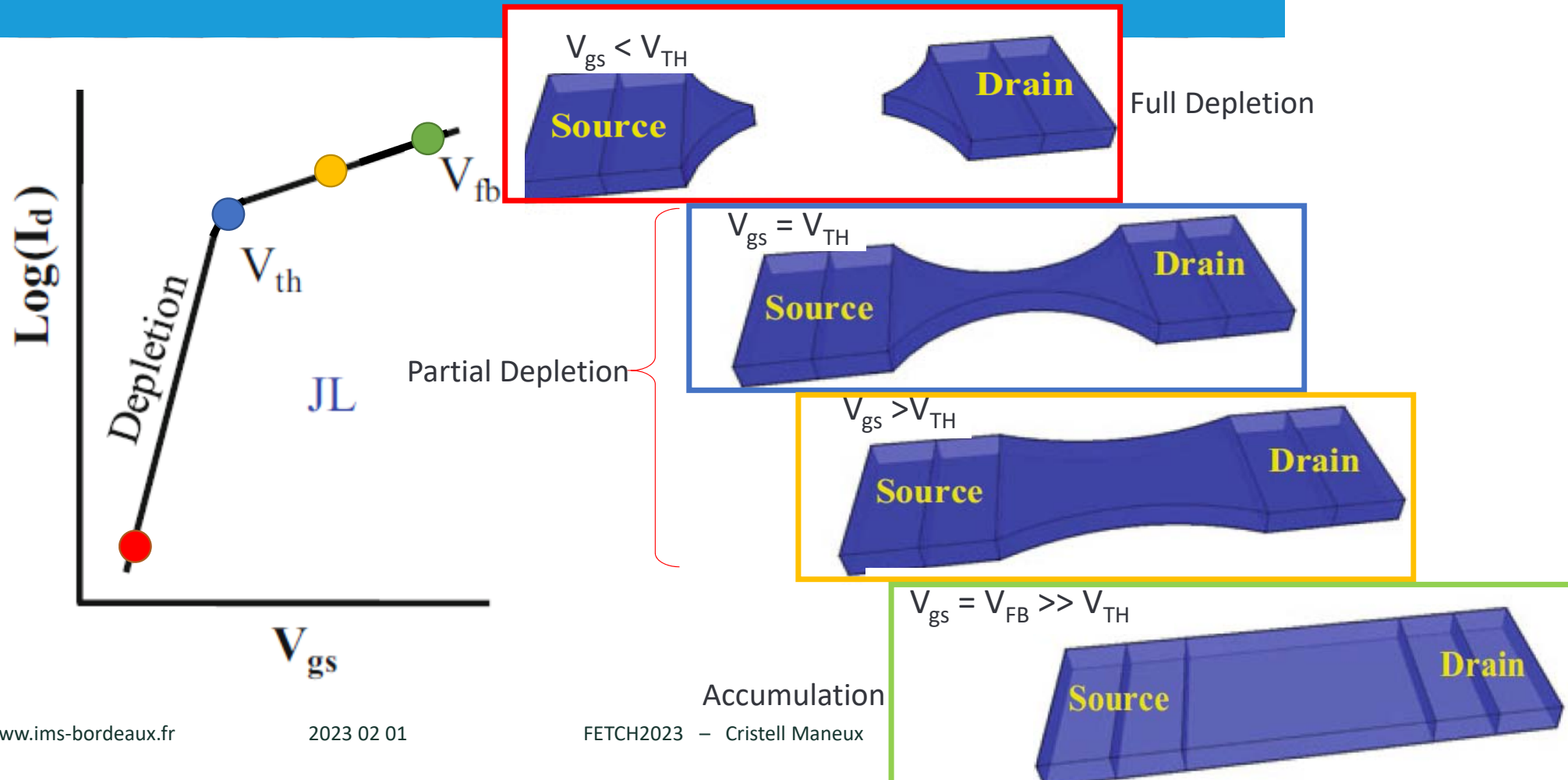


Junctionless FET

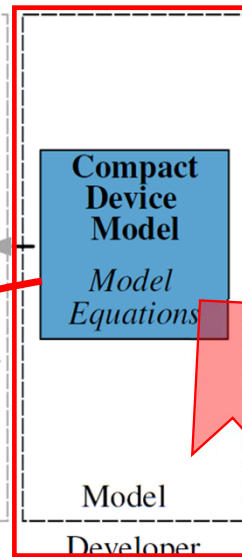
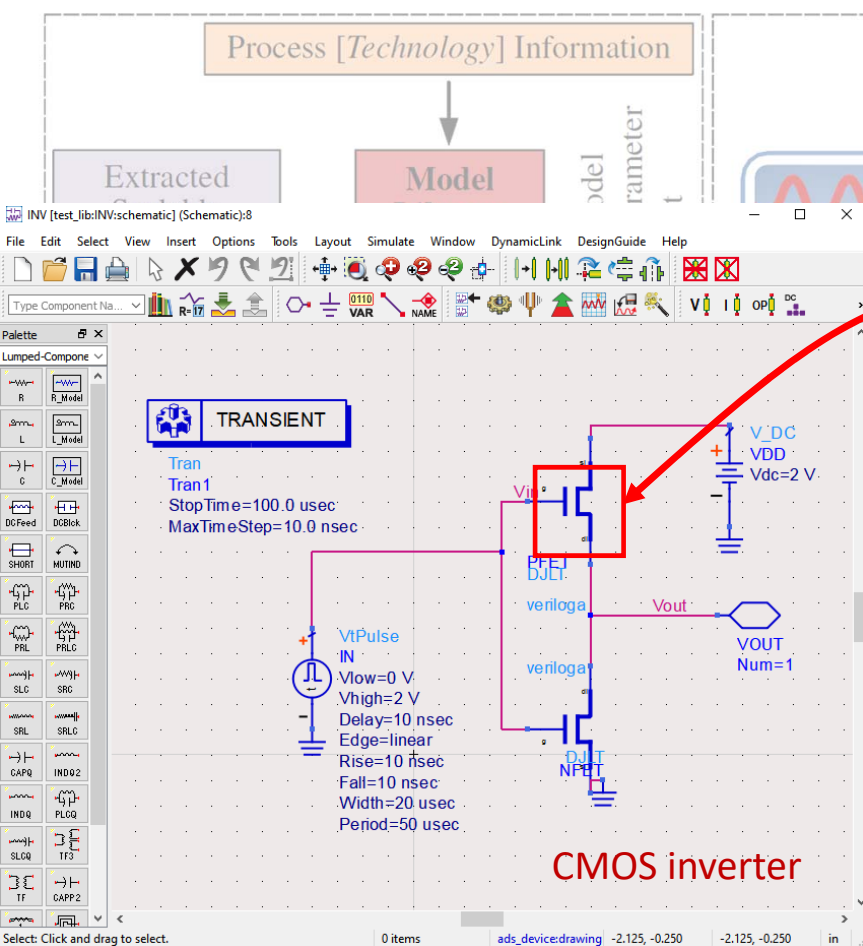


Classical FET

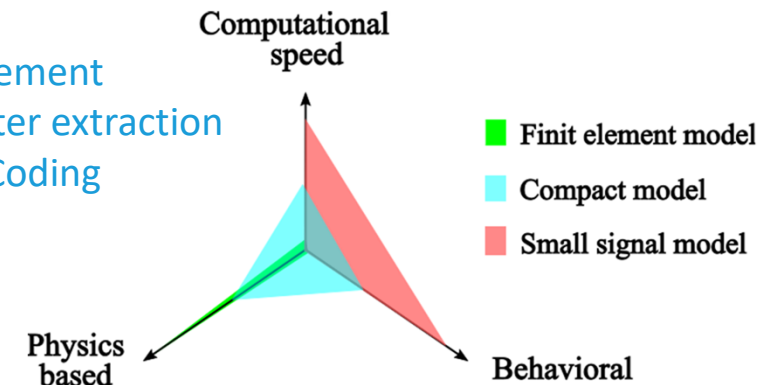
Junctionless FETs: Operation



Compact model: A design tool



- Measurement
- Parameter extraction
- Model Coding



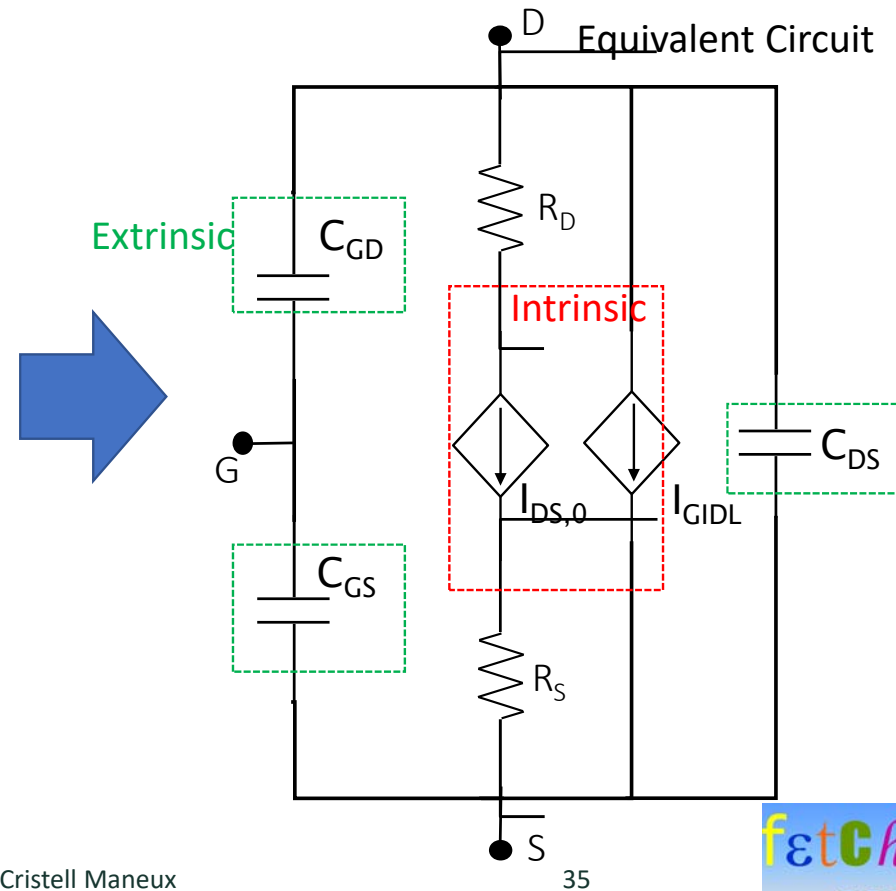
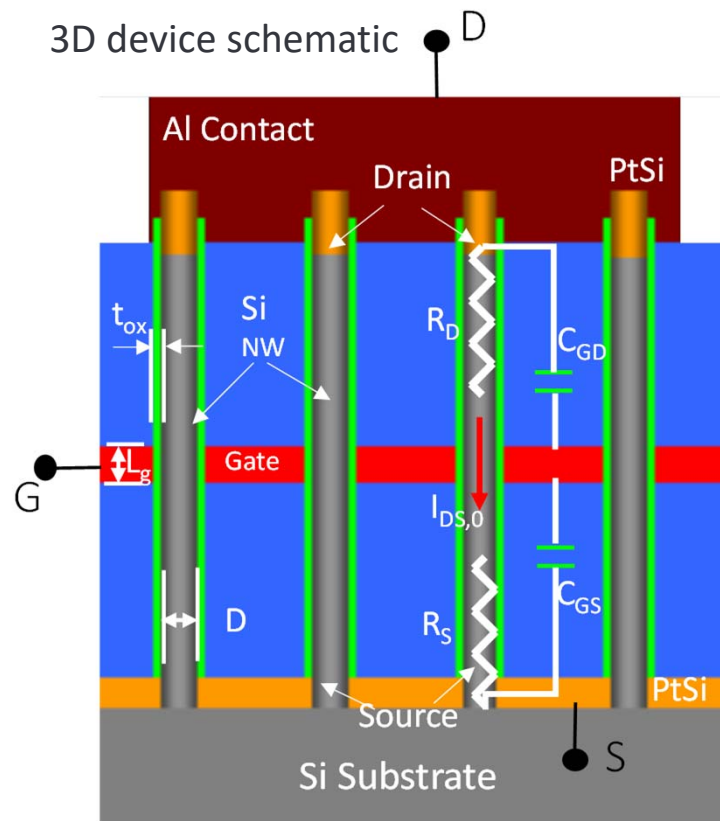
Verilog-A Code

```
//current at source (suffix 0) and drain (suffix L)
IDP0=pow(QDP0,2)/(2*eta2*Cox*uT)+QDP0;
IDPL=pow(QDPL,2)/(2*eta2*Cox*uT)+QDPL;
IC0=pow(QC0,2)/(2*eta2*Cc*uT)+2*QC0;
ICL=pow(QCL,2)/(2*eta2*Cc*uT)+2*QCL;

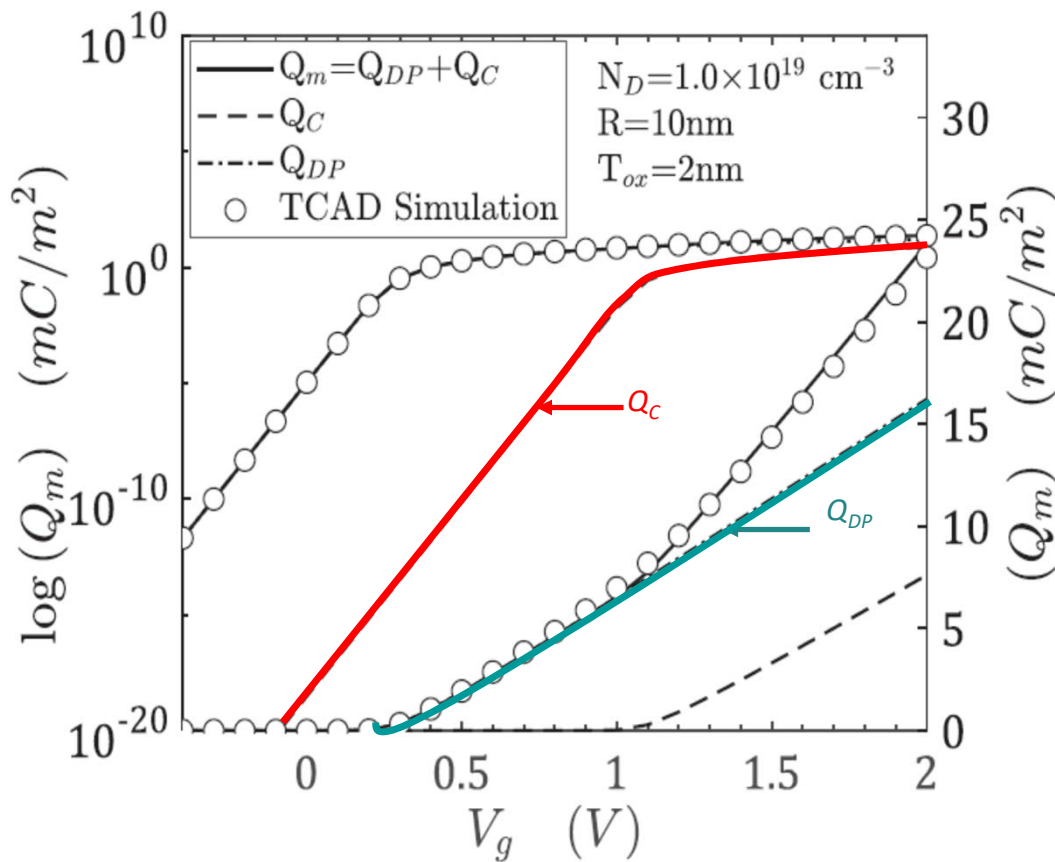
I0=pow(QDP0,2)/(2*eta2*Cox*uT)+QDP0+pow(QC0,2)/(2*eta2*Cc*uT)+2*QC0;
IL=pow(QDPL,2)/(2*eta2*Cox*uT)+QDPL+pow(QCL,2)/(2*eta2*Cc*uT)+2*QCL;

//short channel effect corrections
VS2=(QDP0+QC0)/Cox+Vmax/(Vmax/Vmin-1);
VDSAT=VS2/(VS2/Vmax+1);
VDEFF=VDSAT-VDSAT*(ln(1+exp(A2*(1-vds/VDSAT)))/ln(1+exp(A2)));
```

Junctionless FET compact model



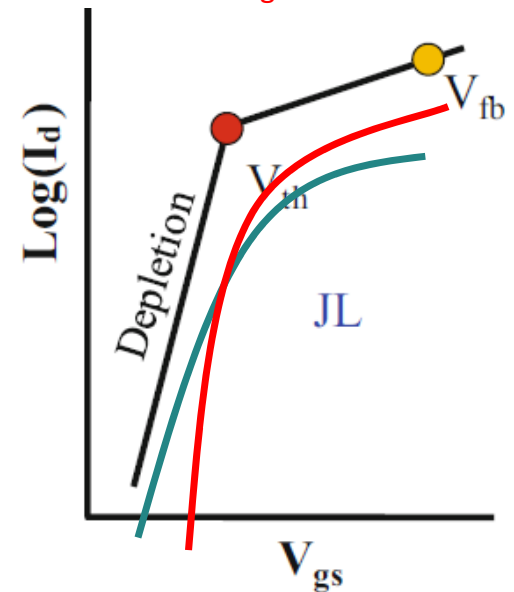
Junctionless FET compact model: I_{DS}



$$I_{DS,0} = \mu_{eff} \frac{2\pi R}{L_{eff}} \left[\int_{Q_{DP0}}^{Q_{DPL}} Q_{DP} \frac{dV}{dQ_{DP}} dQ_{DP} + \int_{Q_{C0}}^{Q_{CL}} Q_C \frac{dV}{dQ_C} dQ_C \right]$$

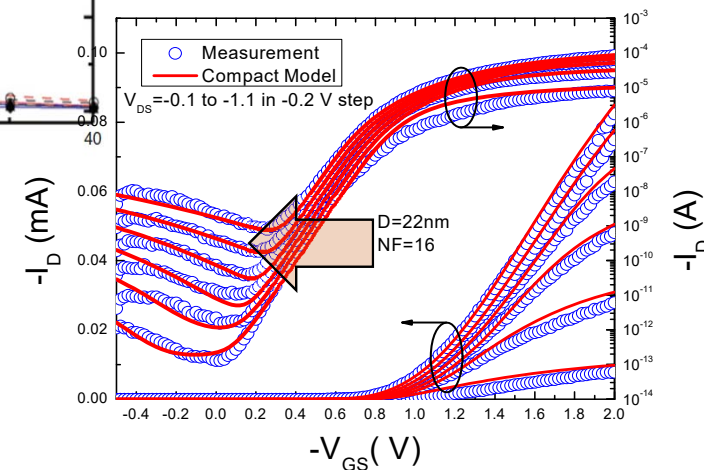
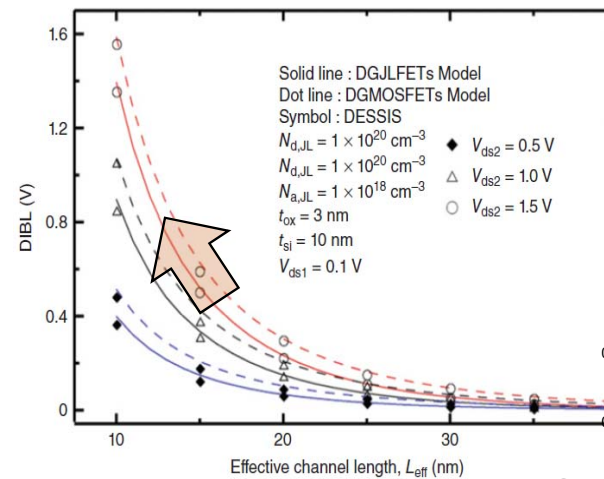
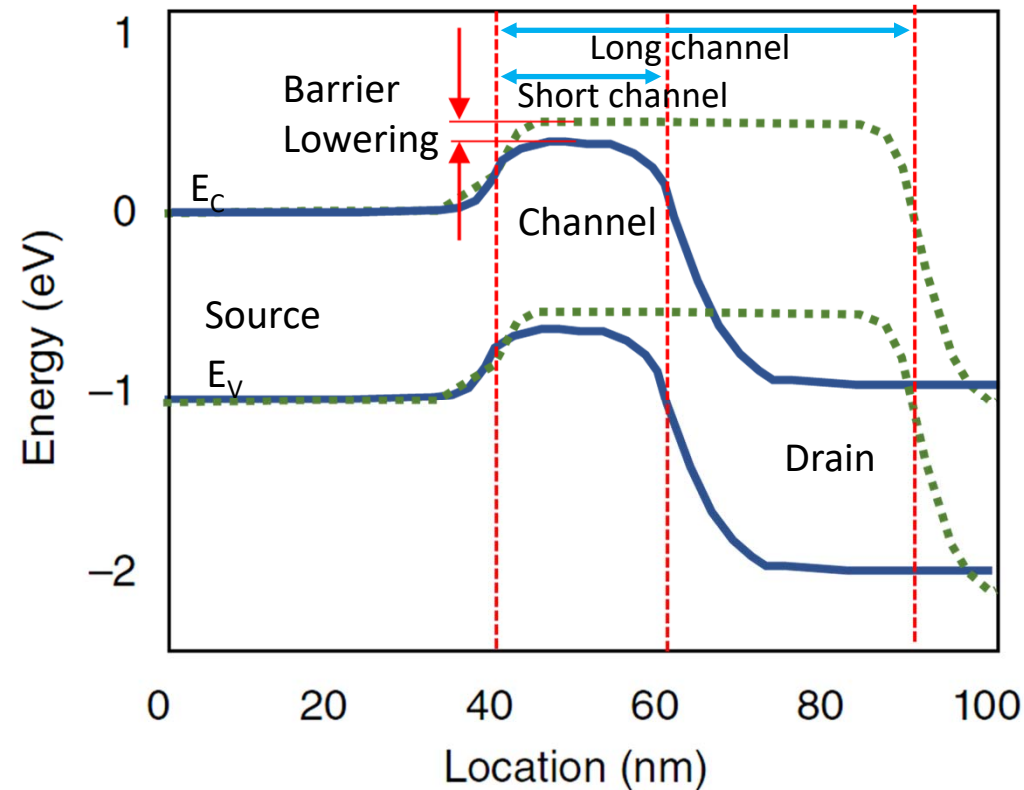
Depletion charge (blue oval) Accumulation charge (red oval)

Circuit diagram of a junctionless FET model. The gate is connected to G, the drain to D, and the source to S. The gate capacitance is C_{GS} , the drain capacitance is C_{GD} , and the source capacitance is C_{DS} . The drain resistance is R_D and the source resistance is R_S . The channel current is $I_{DS,0}$ and the gate-induced drain leakage current is I_{GIDL} .

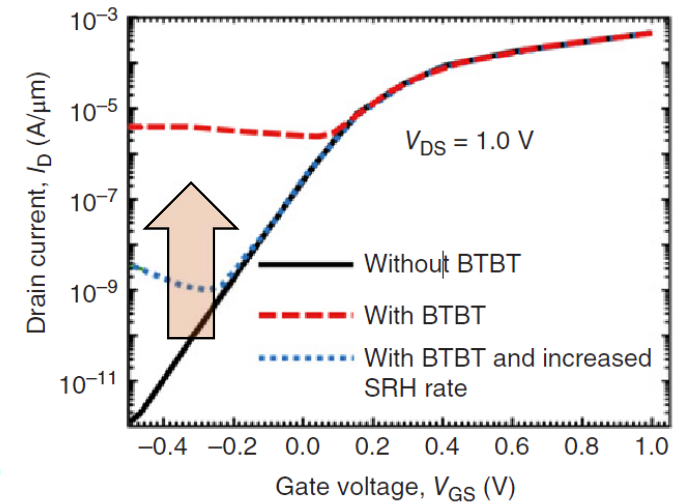
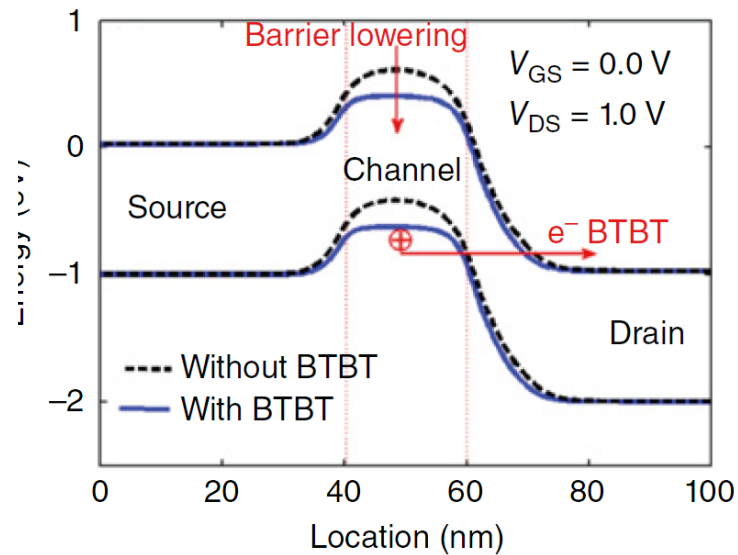
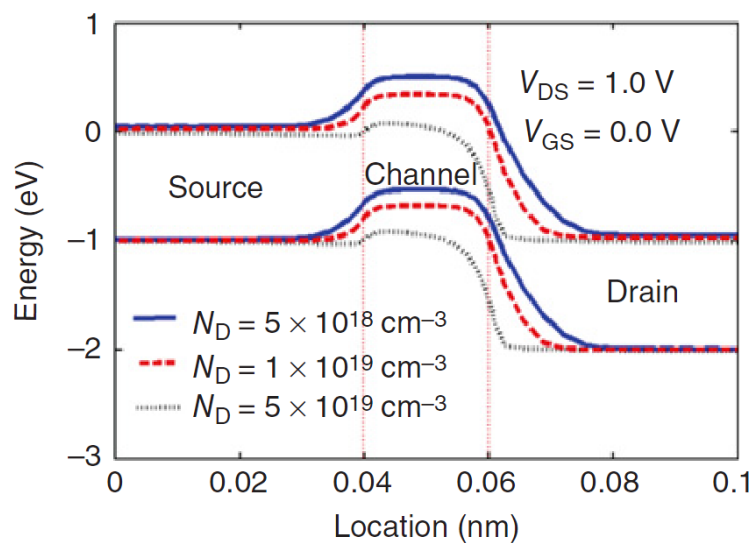


Short Channel Effects: DIBL

Drain-induced barrier lowering



Short Channel Effects: BTBT/GIDL



Gate-induced drain leakage through band-to-band tunneling (BTBT)

Modeling the DC Characteristics

Geometrical parameters

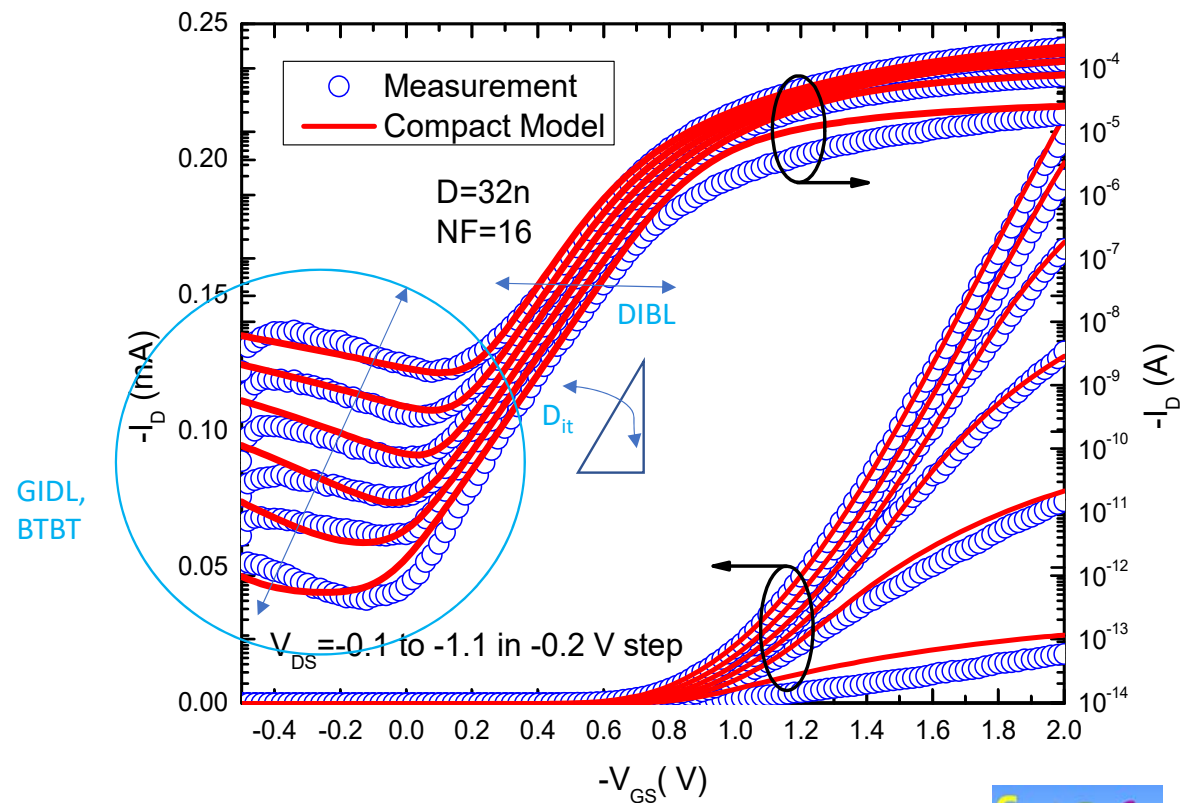
- Diameter (D)
- Number of nanowires in parallel (NF)
- Oxide thickness
- Gate work-function
- NW doping

Modeling the Current

- Mobility
- Carrier saturation velocity
- D/S access resistances

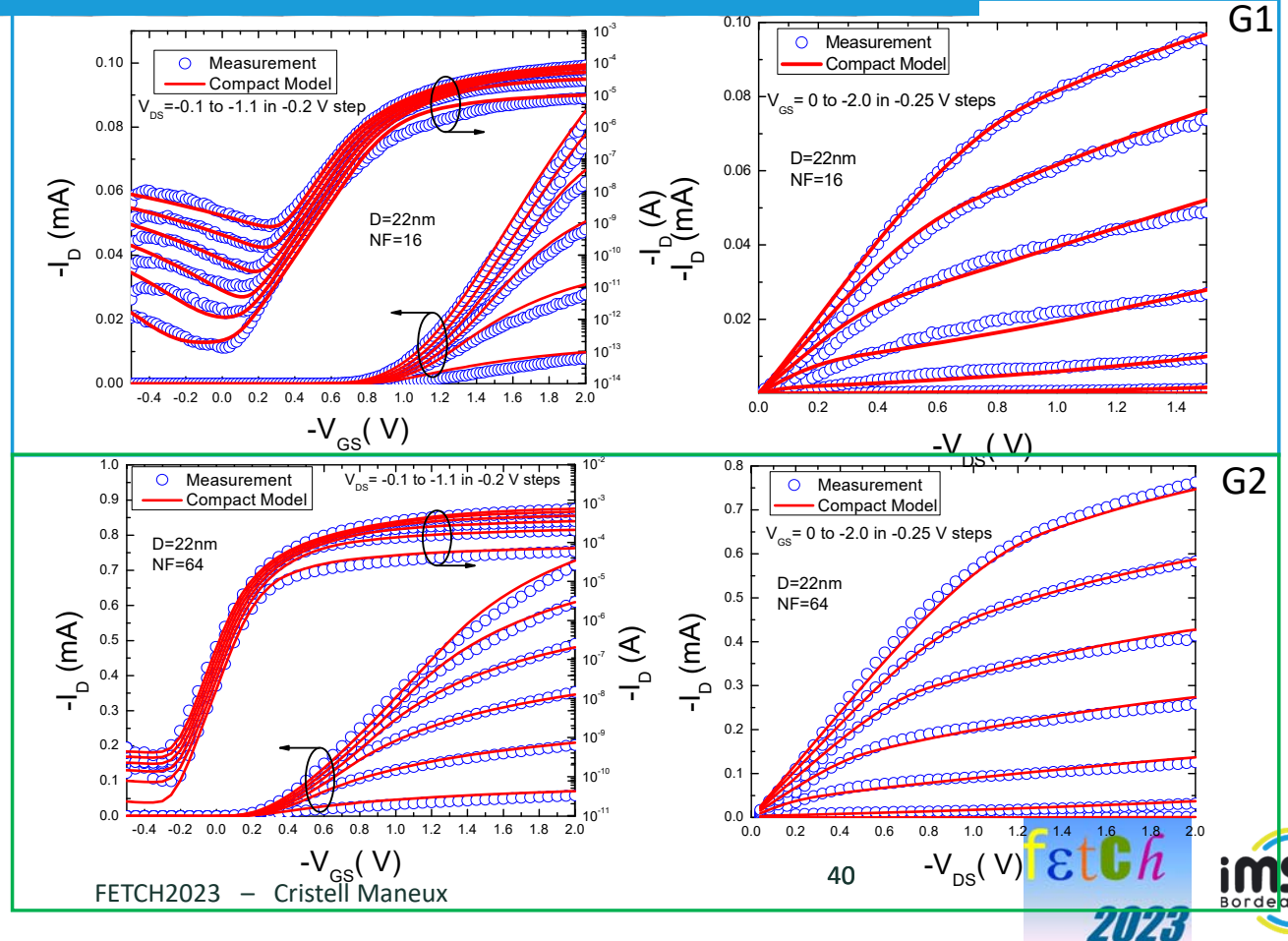
Modeling the Subthreshold:

- Interface Traps
- Drain-induced barrier lowering
- Band-to-band tunneling



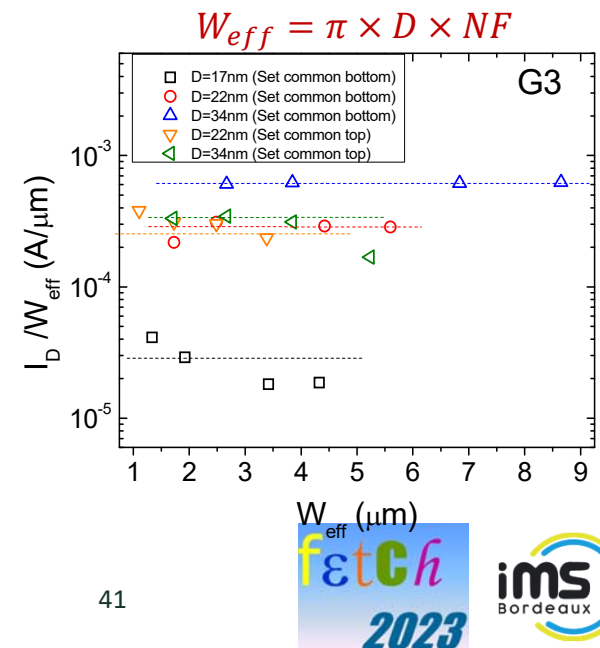
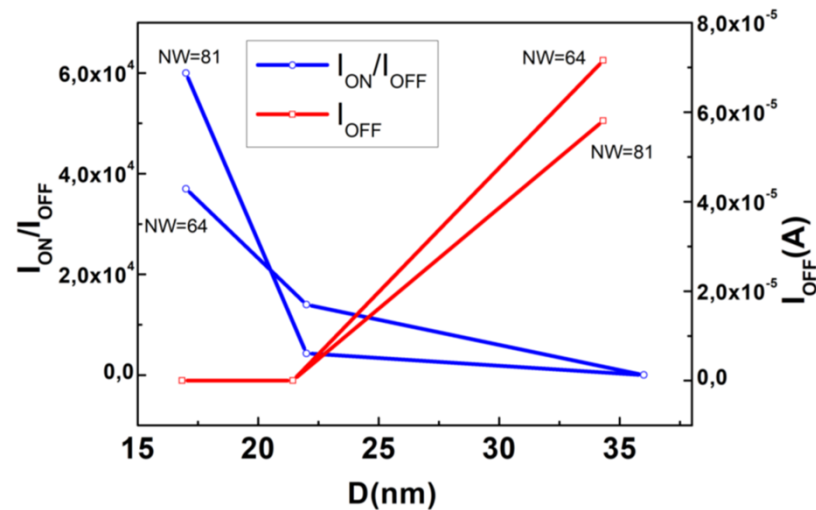
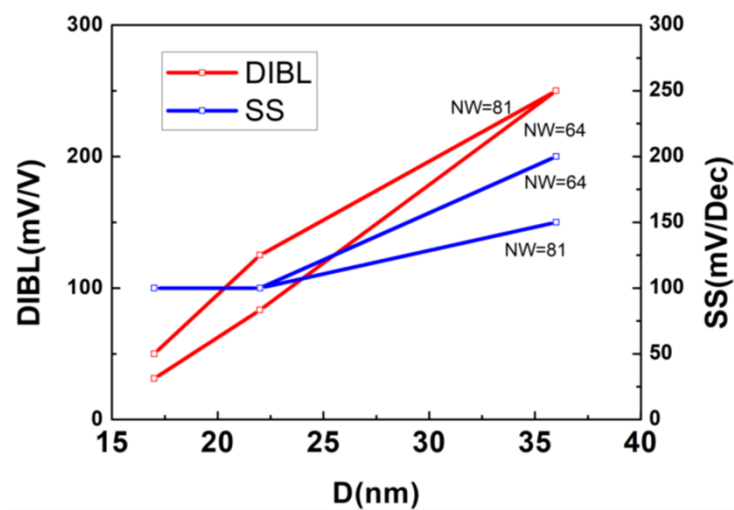
DC Characteristics: I_D - V_G and I_D - V_D

- Two generations of JLNTs
- Improved performance
 - Better subthreshold slope
 - Lower leakage



Scaling

- Scalability is crucial for **predictive modeling**
- Increasing NW diameter degrades DIBL, SS, I_{off} due to loss of electrostatic control
- Optimal design criterion: **smaller diameter** for better Immunity against SCE, **increase of NWs in parallel** for higher drive current and better scaling

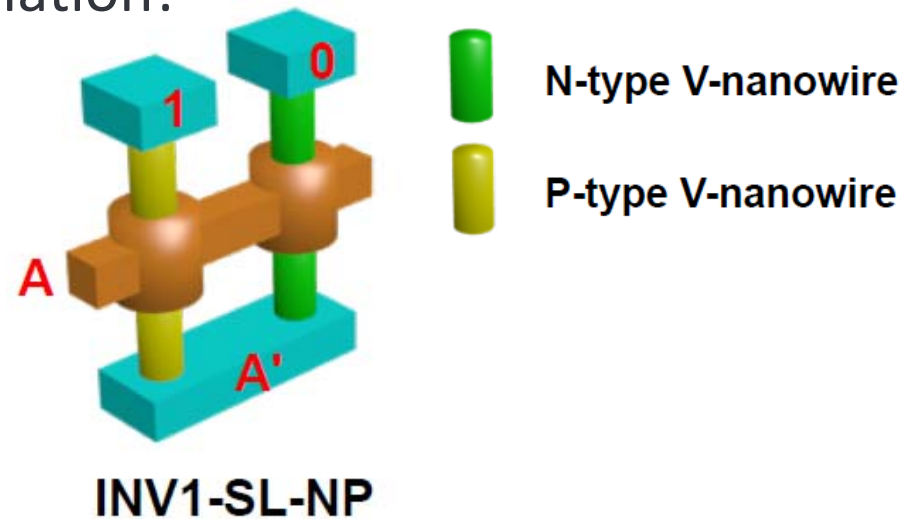


Improvements to Compact Model

- Scaling rules to be improved for **predictive** modeling and circuit design **extrapolation**
- Trap dynamics needs to be studied using **pulse measurements** and modeled for improving drain current behavior at different temperatures
- Electro-thermal effects/**self-heating** and temperature dependence of model parameters (current, threshold voltage , etc.) to be incorporated using an equivalent thermal network.

3D Logic circuits with the VNWFETs

→ Compact model ready for circuit simulation!



Outline


- Contexts
- Characterization challenges
- Modelling challenges
- Circuit design challenges

Towards 3D NN

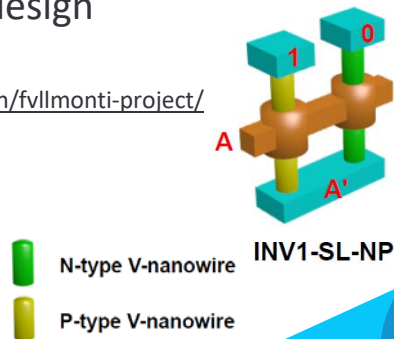
- Dedicated library of 3D logic cells leveraging VNWFET devices
- Versatile VNWFET-logic cell based 3D neural network compute cube (N2C2) for NN-based architecture design
- Scalable and versatile 3D architectural model leveraging reconfigurable 3D interconnect framework
- 3D NN Suitable for exploration of hw/sw co-design

 <https://www.linkedin.com/in/fvllmonti-project/>

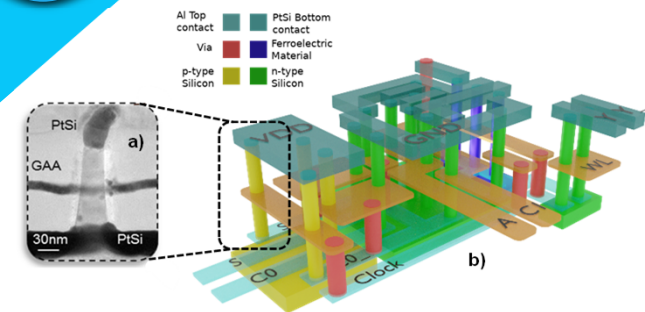
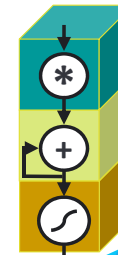
 <https://fvllmonti.eu/>

 Grant n°101016776

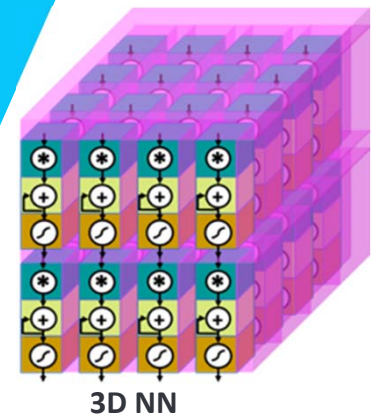
FETCH2022 O'Connor



N²C² concept



Compact and low EDP 1-bit adder



Thanks for listening

