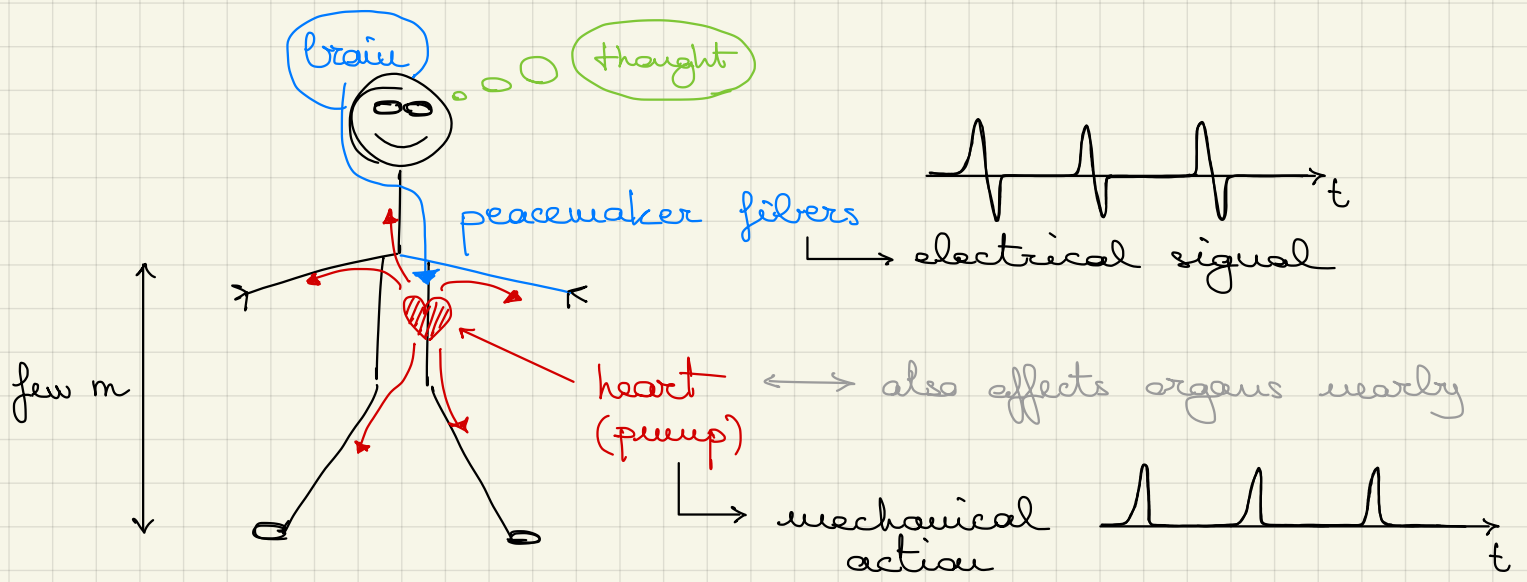


# Computational Modelling in Electronics and Biomechanics

Andrea

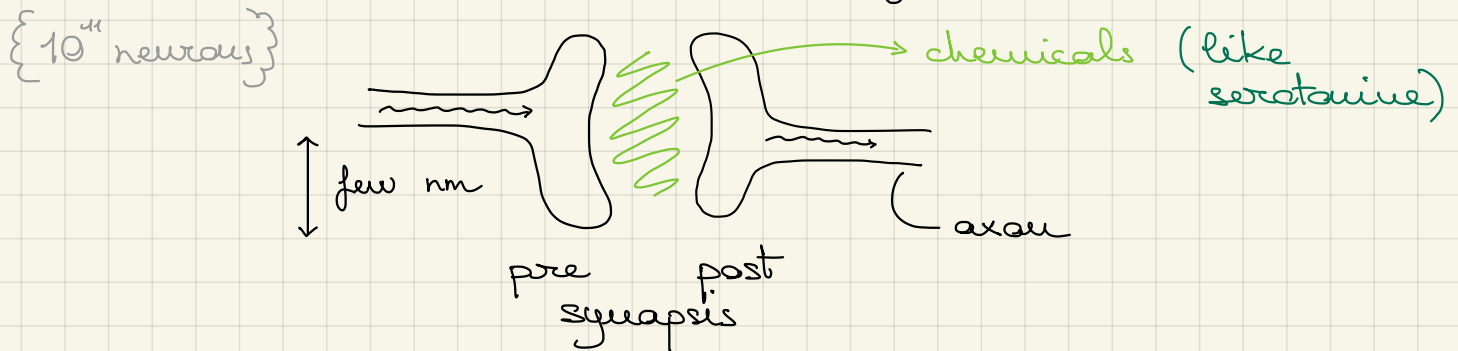
Bertazzoni

A.A. 2020/21



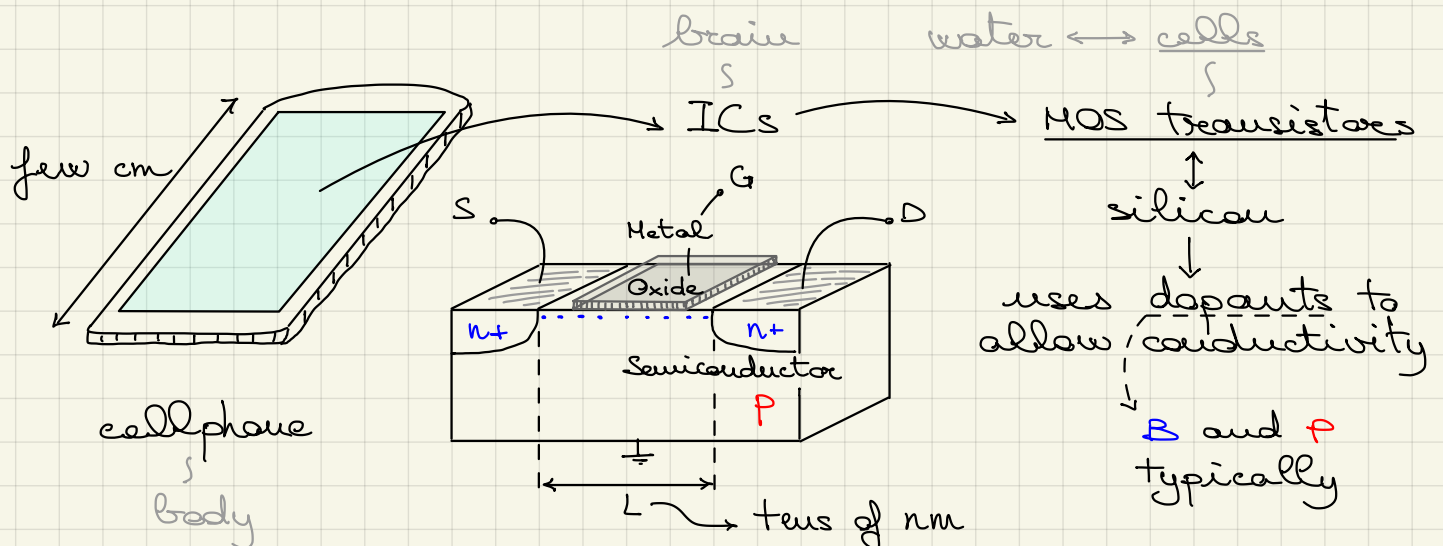
## Electro-mechanical coupling in the human body

### Autonomous nervous system



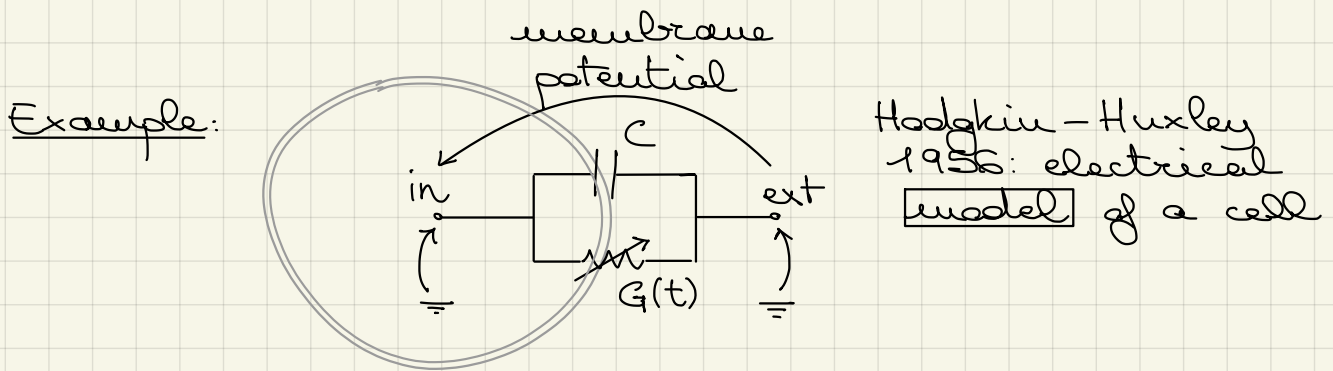
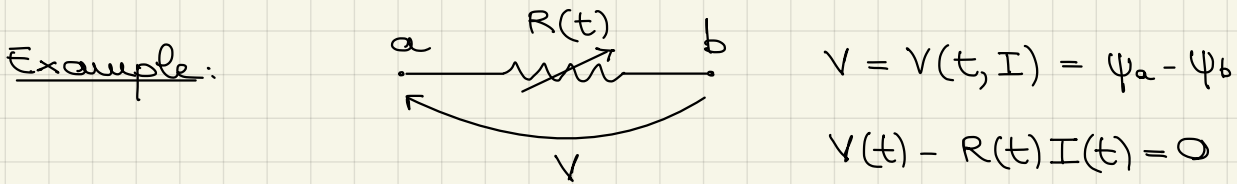
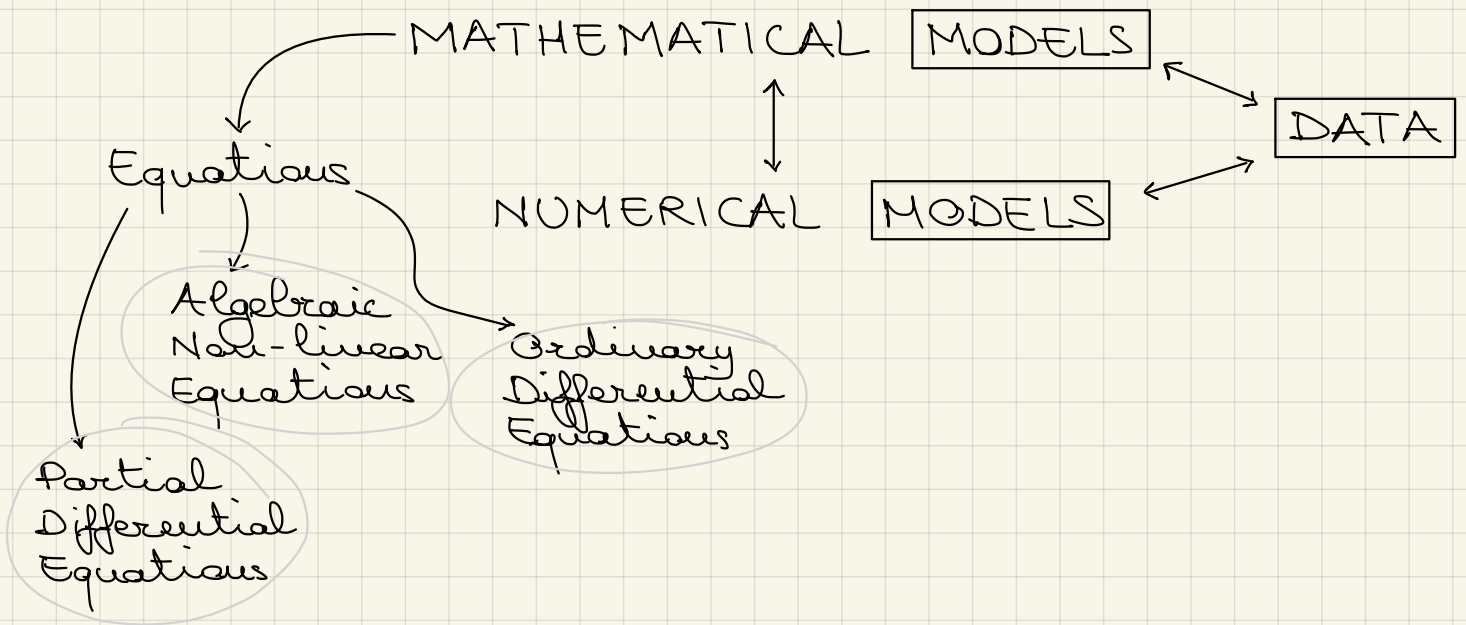
- Chemistry plays a role in human thought and behaviour
- **Blood** and other fluids are pumped through veins and other channels

## Electro-mechanical-chemical-fluidic system





# Numerical Analysis



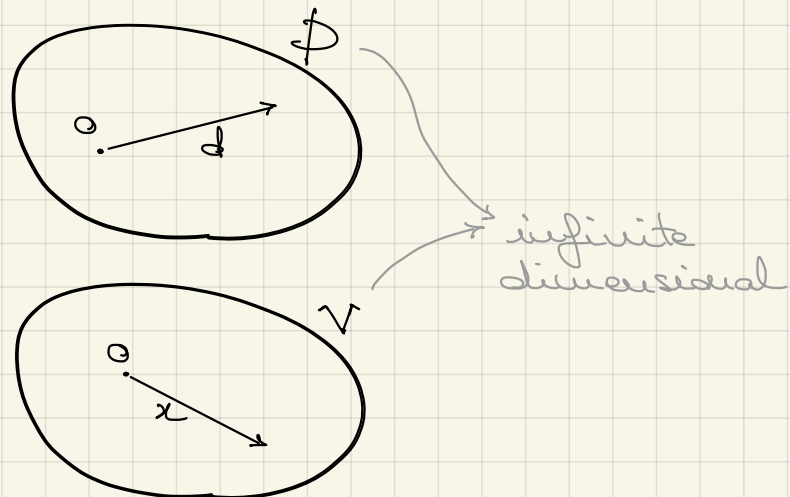
What is a cell model there?  $F(x, d) = 0$

$d$  is a set of data

$\mathcal{D}$  is the set of admissible values for  $d$

$x$  is the unknown

$\mathcal{V}$  is the set of admissible values for  $x$



$$F: (x, d) \longrightarrow y = F(x, d) \iff F: V \times \mathcal{D} \longrightarrow Y$$

Sometimes the closed form of  $F(x, d) = 0$  is not given or is hard to handle.

↳ need of replacing the mathematical (exact) model with the numerical (approximate) model

$$F_h(x_h, d_h) = 0 \quad h > 0$$

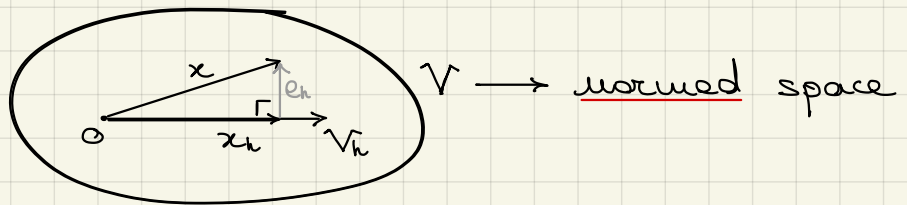
$$e_h := x - x_h \quad \text{error}$$

$$\lim_{h \rightarrow 0} e_h = 0 \iff \text{convergence}$$

$$d_h \in \mathcal{D}_h \subset \mathcal{D} \quad x_h \in V_h \subset V$$

↳ finite dimensional

$F_h$  is a family of functions whose parent function is  $F$ .



$$\phi \in V \quad \text{norm} \leftarrow \|\phi\|_V \geq 0 \quad \|\phi\|_V = 0 \iff \phi = 0$$

$$\left[ \|e_h\|_V \leq c \cdot h^p \right] \quad \underline{p > 0} \quad \text{order of convergence}$$

↳  $c > 0$  indep. of  $h$

$$E_h = c h^p$$

$p = 1$	$h = h_0$	→	$E_0 = c h_0$
linearly converging	$h = \frac{h_0}{2}$	→	$E_1 = c \frac{h_0}{2} = \frac{E_0}{2}$
	$h = \frac{h_0}{4}$	→	$E_2 = c \frac{h_0}{4} = \frac{E_0}{4}$
	⋮	⋮	⋮
	$h = \frac{h_0}{2^n}$	→	$E_n = \frac{E_0}{2^n}$

To obtain  $E = 10^{-6}$  we would need  $h = \frac{10^{-6}}{c}$

An higher convergence order allows to obtain a smaller error with larger discretization step  $h$ .

However higher order methods do not always grant stability in the solution.

Example: test equation

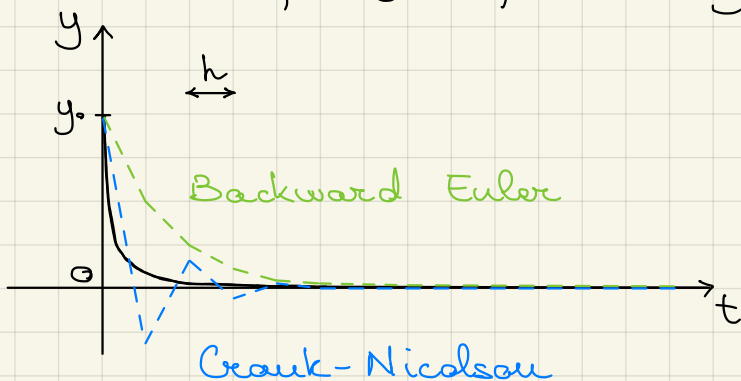
$$(P) \begin{cases} y'(t) = -\lambda y(t) & t > 0 \quad (\lambda > 0) \\ y(t_0) = y_0 \end{cases}$$

math model

dissipative term

initial energy

$$d = \{t_0, y_0, \lambda\} \quad x = y(t) = y_0 e^{-\lambda t}$$



We expect:

$$\lim_{t \rightarrow \infty} y(t) = 0 \quad \checkmark \checkmark$$

$$y(t) \geq 0 \quad \forall t > 0 \quad \checkmark \times$$

Crank-Nicolson is a second order method, hence convergence is faster, however it yields unreliable results for too large values of  $h$ .

Backward Euler is a first order method, hence convergence is slower, but it is always reliable (in this example) for any value of  $h$ .

Convergence  $\Leftrightarrow$  Consistency + Stability

Lax - Richtmyer "Equivalence" theorem

Consistency is granted if a reduction of  $h$  causes a reduction of the residual (= error)

$$\lim_{h \rightarrow 0} e_n = 0$$

Stability is granted if a small variation in the data set causes a small variation in the solution

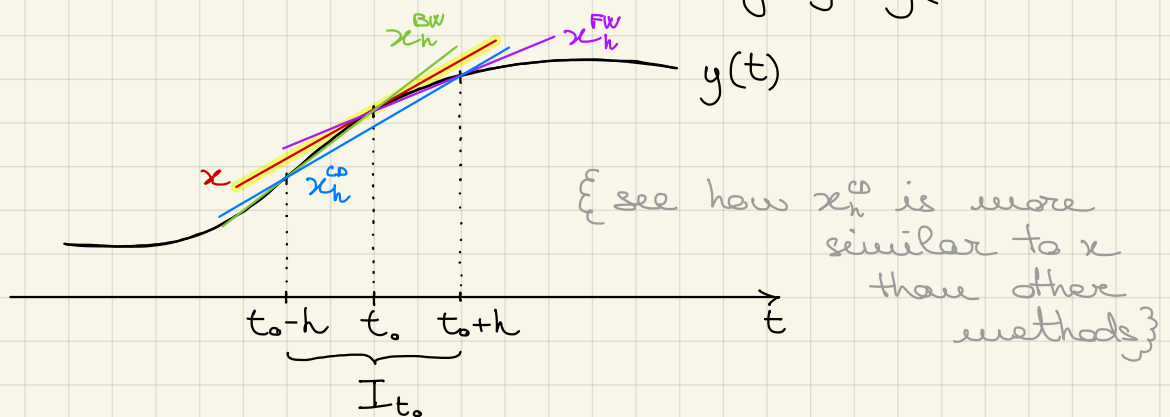
$$\begin{array}{ccc} d & \longrightarrow & d + \delta d \\ \updownarrow & & \updownarrow \\ x = x(d) & & x + \delta x \end{array}$$

If  $\|\delta d\|_p \leq \eta$  with  $\eta$  small there exists  $K = K(d)$  so that

$$\begin{aligned} \|\delta x\|_q &\leq K(d) \|\delta d\|_p \\ &\leq K(d) \eta \end{aligned}$$

↓  
**condition number**

Example: evaluate the derivative of  $y = y(t)$  at  $t = t_0$



Assumption:  $y \in C^2(I_{t_0})$ .

The regularity of the function is of paramount importance when choosing the numerical method.

$$F(x, d) = 0 \quad x = y'(t_0) \in \mathbb{R} = V \quad d = I_{t_0}$$

Taylor's series expansion:  $y(t_0+h) = y(t_0) + \overset{x}{y'(t_0)}h + \frac{y''(\xi)}{2}h^2$   
(to second order)

where  $\xi \in [t_0, t_0+h]$

and  $|y''(\xi)| < M$  since  $y \in C^2(I_{t_0})$

$$\Rightarrow \frac{y(t_0+h) - y(t_0)}{h} - x - \underbrace{\frac{y''(\xi)}{2}h^2}_{\text{error}} = 0 \quad \Leftrightarrow \quad F(x, d) = 0$$

Mathematical model

$$\Rightarrow \frac{y(t_0+h) - y(t_0)}{h} - x_{\underline{n}} = 0 \quad \Leftrightarrow \quad F_n(x_n, d_n) = 0$$

Numerical model

"Forward" finite difference formula:

$$\left[ x_h^{FW} = \frac{y(t_0+h) - y(t_0)}{h} \right]$$

"Backward" finite difference formula:

$$\left[ x_h^{BW} = \frac{y(t_0) - y(t_0-h)}{h} \right]$$

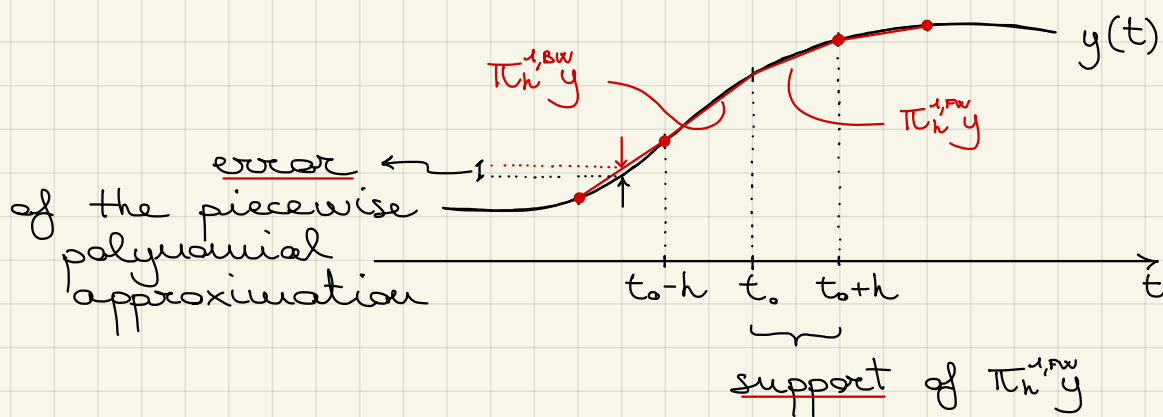
$$|e_n| = |y'(t_0) - x_n| \leq \underbrace{c}_{\frac{M}{2}} h \longrightarrow \text{first order methods}$$

"Centered" finite difference formula:

$$\left[ x_h^{CD} = \frac{y(t_0+h) - y(t_0-h)}{2h} \right]$$
$$= \frac{1}{2} [x_h^{FW} + x_h^{BW}]$$

$$|e_n| = c h^2 \longrightarrow \text{second order method (requires stricter conditions on } y)$$

What we are doing when using these finite difference methods is actually a piecewise polynomial interpolation



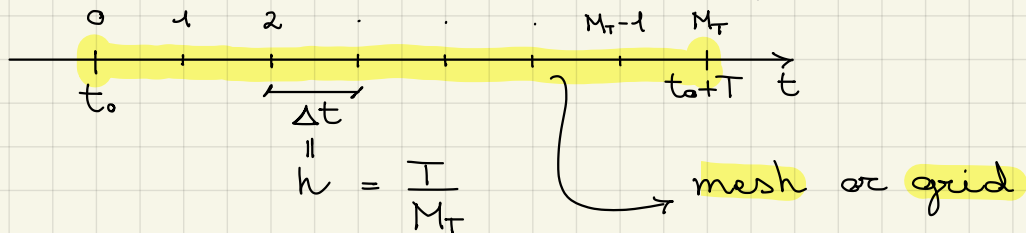
$$x_h^{FW} = \left. \frac{d}{dt} \pi_h^{1,FW} y(t) \right|_{t=t_0}$$

$$x_h^{BW} = \left. \frac{d}{dt} \pi_h^{1,BW} y(t) \right|_{t=t_0}$$

Application of this example: Cauchy Problem

$$\begin{cases} y'(t) = f(t, y(t)) & t \in I_T = (t_0, t_0+T) \\ y(t_0) = y_0 \end{cases}$$

if  $f(t, y(t)) = -\lambda y(t)$   
the problem becomes the test equation



$\begin{cases} t_k = t_0 + k\Delta t \\ k = 0, 1, \dots, M_T \end{cases}$  we are just going to evaluate the solution of the problem at these nodes (the value in-between is not known).

$u_k :=$  approximate solution       $y_k := y(t_k) =$  real solution

$$u_0 = y_0$$

only certainty we have

$$t_0 \longrightarrow t_1 = t_0 + h \longrightarrow t_2 = t_1 + h \longrightarrow t_3 = t_2 + h \longrightarrow \dots$$

$$u_0 = y_0 \qquad u_1 \qquad u_2 \qquad u_3$$

"Forward Euler" method

$$\frac{u_1 - u_0}{h} = f(t_0, u_0) \longrightarrow \begin{cases} u_{k+1} = u_k + h f(t_k, u_k) & k = 0, 1, \dots, M_T-1 \\ u_0 = y_0 \end{cases}$$

$$\frac{u_1 - u_0}{h} = f(t_1, u_1) \longrightarrow \begin{cases} u_{k+1} = u_k + h f(t_{k+1}, \underline{u_{k+1}}) & k = 0, 1, \dots, M_T-1 \\ u_0 = y_0 \end{cases}$$

\*

"Backward Euler" method

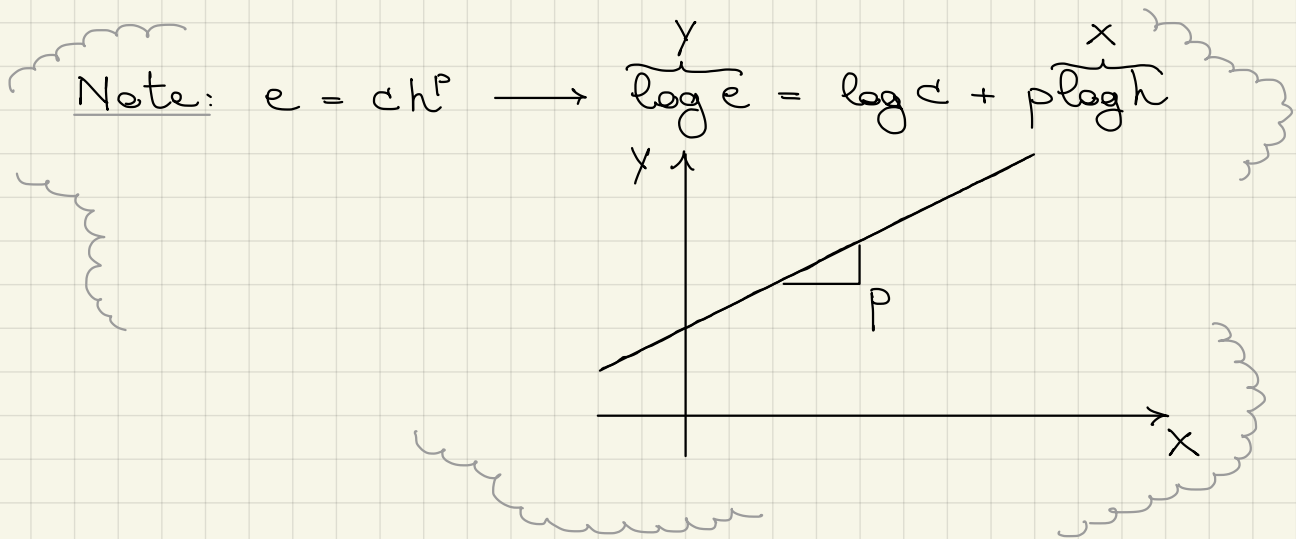
\* Algorithm where the solution is not explicitly computable; in other words, it has the form of an algebraic differential equation (if  $f$  is non-linear with respect to  $y$ ).

## "Crank-Nicolson" method

$$\begin{cases} u_{k+1} = u_k + \frac{h}{2} [f(t_k, u_k) + f(t_{k+1}, u_{k+1})] & k = 0, 1, \dots, M_T - 1 \\ u_0 = y_0 \end{cases}$$

$|e_k| = |y_k - u_k| \leq ch$  for Forward/Backward Euler (1st order)

$|e_k| = ch^2$  for Crank-Nicolson (2nd order)



We already saw that a drawback of the Crank-Nicolson method was that it oscillates for too large  $h$ .

→ Crank-Nicolson is not a monotone method

↕  
positivity preserving

→ Trade-off between convergence order and monotonicity

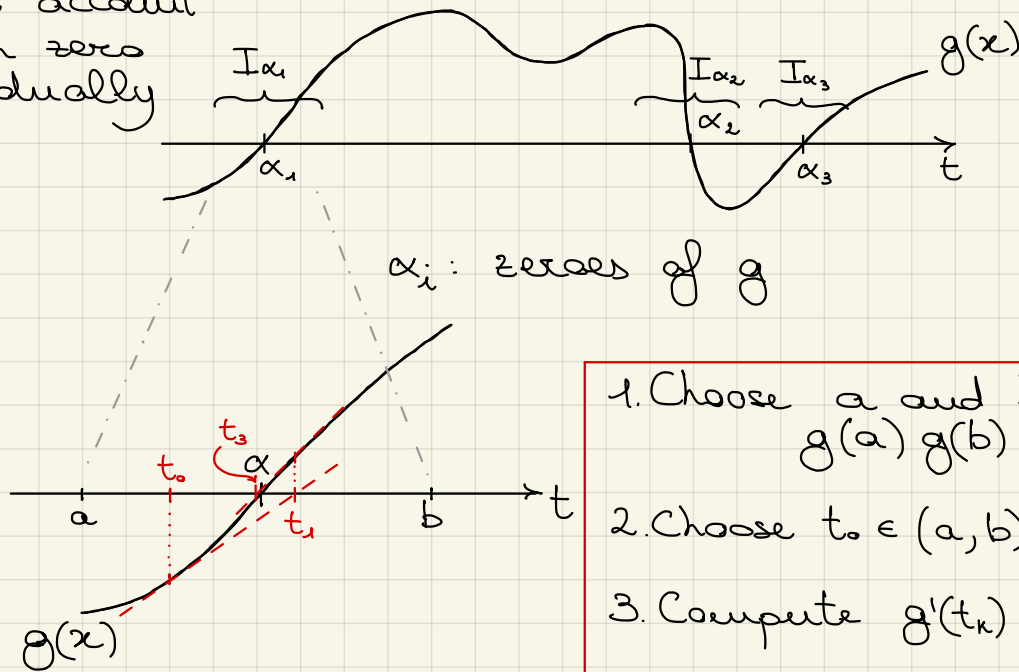
Backward Euler:  $u_{k+1} = u_k + h f(t_{k+1}, u_{k+1})$

$$\boxed{x - u_k - h f(t_{k+1}, x) = 0}$$

Algebraic differential equation

$$\rightarrow \boxed{g(x) = 0}$$

Need to account for each zero individually



1. Choose  $a$  and  $b$  so that  $g(a)g(b) < 0$
2. Choose  $t_0 \in (a, b)$
3. Compute  $g'(t_k)$  and  $g(t_k)$
4. Find intersect of tangent at  $g(t_k)$  with  $t$  axis and call it  $t_{k+1}$
5. Iterate from point 3 until  $g(t_k) \approx 0$

second order method

"Newton's" method  
(method of tangents)

Even if the method does not reach exactly zero it will anyways reach machine precision

smallest number displayed by machine's bits architecture (e.g. 64 bits  $\sim 10^{-16}$ )

$\theta$ -method: a general form for numerical methods for differential equations

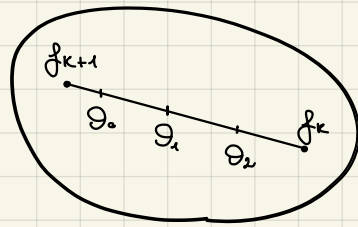
$$\theta \in [0, 1]$$

$$\begin{cases} y'(t) = f(t, y(t)) & t \in I_T = (t_0, t_0 + T) \\ y(t_0) = y_0 \end{cases}$$

$$\begin{cases} u_{k+1} = u_k + h \left[ \theta \underbrace{f(t_{k+1}, u_{k+1})}_{f_{k+1}} + (1-\theta) \underbrace{f(t_k, u_k)}_{f_k} \right] & k = 0, 1, \dots, M_T - 1 \\ u_0 = y_0 \end{cases}$$



$[\vartheta f_{k+1} + (1-\vartheta)f_k]$  is a family of functions:



$\vartheta = 0 \longrightarrow$  Forward Euler

$\vartheta = \frac{1}{2} \longrightarrow$  Crank-Nicolson

$\vartheta = 1 \longrightarrow$  Backward Euler

Theorem: There exists a positive constant  $c$  INDEPENDENT of  $h$  such that

$$|e_k| \leq c h^{p(\vartheta)}$$

$$\text{where } p(\vartheta) = \begin{cases} 1 & \vartheta \neq \frac{1}{2} \\ 2 & \vartheta = \frac{1}{2} \end{cases}$$

Absolute stability  $\neq$  stability

$$(P) \begin{cases} y'(t) = -\lambda y(t), & t > t_0 \\ y(t_0) = y_0 \end{cases} \implies \begin{matrix} T = +\infty \\ f(t, y(t)) = -\lambda y(t) \\ \lambda > 0 \quad (\lambda \in \mathbb{R}) \end{matrix}$$

Use  $\vartheta$ -method to find the solution of the problem

$$y(t) = y_0 e^{-\lambda(t-t_0)}$$

We know that the mathematical solution is asymptotically stable since  $\lambda > 0$  (Lyapunov stable).

$$\lim_{t \rightarrow +\infty} y(t) = 0$$

What can we say about the numerical solution?

$$\lim_{n \rightarrow +\infty} u_n = 0 \iff \text{absolute stability}$$

$$f(y) = -\lambda y \quad f_n = -\lambda u_n \quad f_{n+1} = -\lambda u_{n+1}$$

$$\theta\text{-method: } \begin{cases} \frac{u_{n+1} - u_n}{h} = -\theta \lambda u_{n+1} - (1-\theta) \lambda u_n & n \geq 0 \\ u_0 = y_0 \end{cases}$$

$$\longrightarrow u_{n+1} \left[ \frac{1}{h} + \theta \lambda \right] = u_n \left[ \frac{1}{h} - (1-\theta) \lambda \right]$$

$$\implies u_{n+1} = u_n \frac{1 - \lambda h (1-\theta)}{1 + \lambda h \theta} \quad n \geq 0$$

$$= u_n \cdot \phi_\theta(\lambda, h)$$

$$u_1 = u_0 \phi_\theta(\lambda, h) \longrightarrow u_2 = u_1 \phi_\theta(\lambda, h) = u_0 [\phi_\theta(\lambda, h)]^2$$

$$\implies u_n = u_0 [\phi_\theta(\lambda, h)]^n \quad n \geq 0$$

$$y_n = y_0 e^{-\lambda n h}$$

To grant absolute stability it must be  $\lim_{n \rightarrow +\infty} u_n = 0$

hence  $|\phi_\theta(\lambda, h)| < 1 \implies -1 < \frac{1 - \lambda h (1-\theta)}{1 + \lambda h \theta} < 1$

1)  $-1 - \lambda h \theta < 1 - \lambda h + \lambda h \theta$       2)  $1 + \lambda h \theta > 1 - \lambda h + \lambda h \theta$

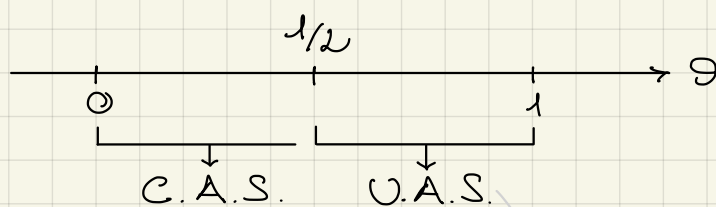
$$h < \frac{2}{\lambda} \frac{1}{1-2\theta} \quad h > 0$$

$$\implies 0 < h < \frac{2}{\lambda} \frac{1}{1-2\theta}$$

can be chosen  $0 \leq \theta \leq 1$

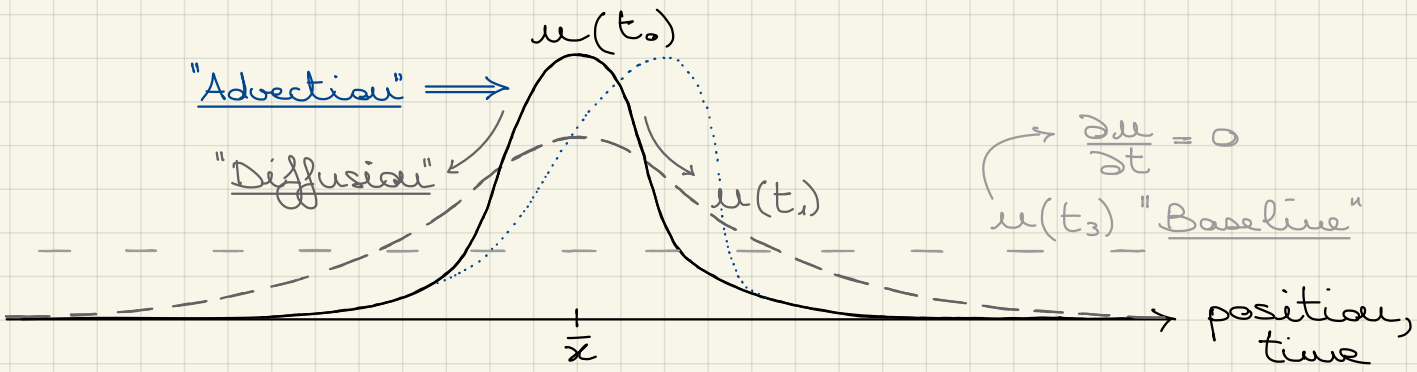
$\theta = 0 \longrightarrow 0 < h < \frac{2}{\lambda}$

$\theta = \frac{1}{2}$   $\longrightarrow$  unconditional absolute stability

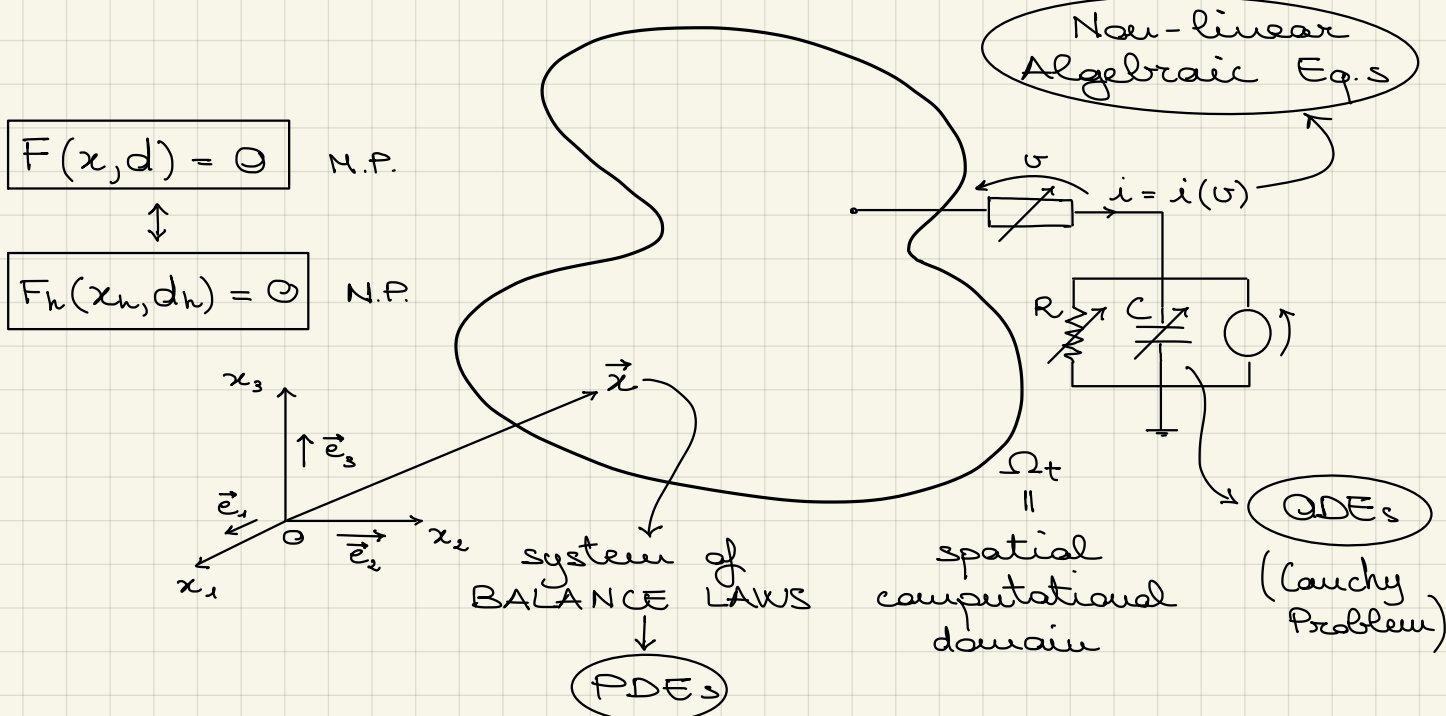


for  $\theta > \frac{1}{2}$  the right-hand inequation doesn't make sense so we only consider the left-hand one

# Physical Mechanisms



- Diffusion → Fick's Law
  - Advection
  - Net production =  $G - R$       $G, R > 0$ 
    - ↑ generation
    - ↙ recombination
- NO mechanical phenomena



$\Omega_t$  should not be considered as a fixed domain but it has to be described by an evolution equation:

$$\Omega_0 = \Omega_t|_{t=0} \leftrightarrow t \in [0, T]$$

For now, we will not consider the time dependency of  $\Omega_t$  and the non-linear differential couplings with the outer world.

$$u = u(\vec{x}, t)$$

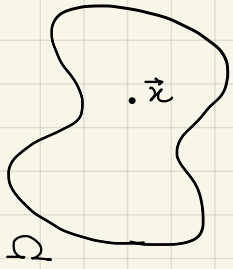
$$\vec{j} = \vec{j}(u)$$

$$\frac{\partial u}{\partial t} + \vec{\nabla} \cdot \vec{j} = \rho$$

$$\vec{x} \in \Omega \subset \mathbb{R}^3$$

$$t \in [0, T]$$

$$\vec{j} = \begin{bmatrix} j_1(\vec{x}, t) \\ j_2(\vec{x}, t) \\ j_3(\vec{x}, t) \end{bmatrix}$$



$$\frac{\partial j_1}{\partial x_1} + \frac{\partial j_2}{\partial x_2} + \frac{\partial j_3}{\partial x_3}$$

time rate of change of  $u$

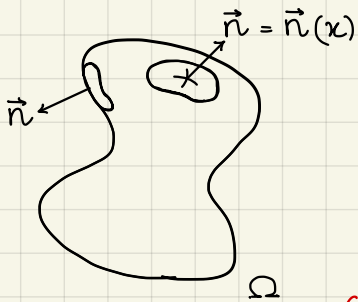
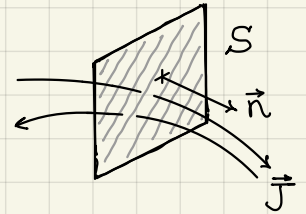
Do mind physical dimensions!

$$[u] = U \quad [\rho] = \frac{U}{S} \quad [j] = U \cdot \frac{m}{S}$$

$\Rightarrow \vec{j}$ : flux density of  $u$

$$\Phi_u^S = \int_S \vec{j} \cdot \vec{n} d\Sigma = \text{flux of } u \text{ across } S$$

$$[\Phi_u] = U \frac{m}{S} \cdot m^2 = U \frac{m^3}{S}$$



$$\int_{\Omega} \vec{\nabla} \cdot \vec{j} d\Omega = \int_{\partial\Omega} \vec{j} \cdot \vec{n} d\Sigma = \Phi_u^{\partial\Omega} = \Phi_u^{\partial\Omega}(t)$$

Stokes theorem

$$\begin{cases} \frac{\partial u}{\partial t} + \vec{\nabla} \cdot \vec{j} = \rho & \text{balance law} \\ \vec{j} = \dots & \text{constitutive law for } \vec{j} \\ \rho = \dots & \text{" " " " } \rho \end{cases}$$

We are still lacking: ① Initial condition

and

② Boundary conditions

$t=0$

①  $u(\vec{x}, 0) = u^i(\vec{x})$  given function of  $\vec{x} \in \Omega$

② On the boundary  $\partial\Omega$ :

$$\begin{array}{ccc} \hookrightarrow u & \hookrightarrow \vec{j} & \hookrightarrow \text{both} \\ \partial\Omega = \partial\Omega_D \cup \partial\Omega_N \cup \partial\Omega_R & & \\ \uparrow & \uparrow & \uparrow \\ \text{Dirichlet} & \text{Neumann} & \text{Robin} \end{array}$$

Dirichlet b.c.:  $u|_{\partial\Omega_D} = \bar{u}_0(\vec{x}, t) \quad \vec{x} \in \partial\Omega_D$

Neumann b.c.:  $\vec{j} \cdot \vec{n}|_{\partial\Omega_N} = \bar{j}_N(\vec{x}, t) \quad \vec{x} \in \partial\Omega_N$

Robin b.c.:  $\bar{\alpha}_R \vec{j} \cdot \vec{n}|_{\partial\Omega_R} = \bar{\alpha}_R(\vec{x}, t) u - \bar{\beta}_R(\vec{x}, t) \quad \vec{x} \in \partial\Omega_R$

We now need to explicit the constitutive laws.

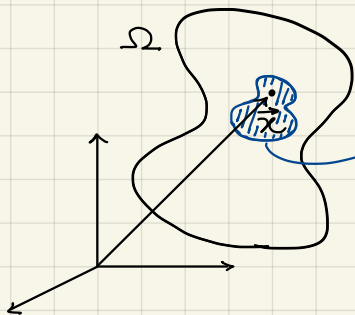
$\vec{j} = \vec{V}u - D\vec{\nabla}u \rightarrow$  advection-diffusion model for the flux density of  $u$

$\vec{j}_a = \vec{V}u \rightarrow$  velocity field  $[\vec{V}] = \frac{m}{s}$

$\vec{j}_d = -D\vec{\nabla}u \rightarrow$  Fick's law

$\rightarrow$  diffusion coefficient  $[D] = \frac{m^2}{s}$

$\Phi = ? \rightarrow$  microscopic model of the system



"material volume"  
 $dV_k$

"The smallest possible volume that retains the properties and physics of the material"



mass density  
 $\rho^m(\vec{x}, t) = \lim_{R \rightarrow 0^+} \frac{m_{dV_k}}{dV_k}$

$\downarrow$   
 $\rho^q(\vec{x}, t) = \lim_{R \rightarrow 0^+} \frac{Q_{dV_k}}{dV_k}$   
charge density

$\rightarrow \Phi = g - \kappa u$  our chosen model

production

$g = g(\vec{x}, t) > 0$

$\kappa = \kappa(\vec{x}, t) > 0$

$\Phi = \Phi(\vec{x}, t)$

consumption

Note: if  $\Phi(\vec{x}_0, t_0) = 0$ , then  $u(\vec{x}_0, t_0) = \frac{g(\vec{x}_0, t_0)}{\kappa(\vec{x}_0, t_0)}$

Under the assumptions that: 1)  $\vec{\nabla} \cdot \vec{j}$  is negligible and 2) the system has reached steady-state ( $\frac{\partial u}{\partial t} = 0$ ) the solution  $u(\vec{x}, t) = u(t) = u(+\infty)$  is exactly that of  $\Phi = 0$ .

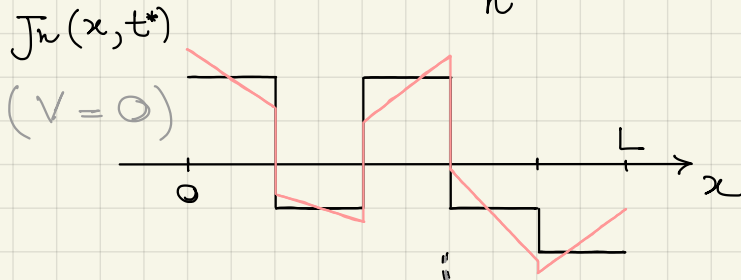
$\leftarrow$  "lumped" model i.e. does not depend on  $x$





Solving for  $u$ :

$u \rightarrow u_n$   
finite element approximation



Derive  $J_n$  from  $u_n$ :

$$J_n = J(u_n) = -D \frac{\partial u_n}{\partial x}$$

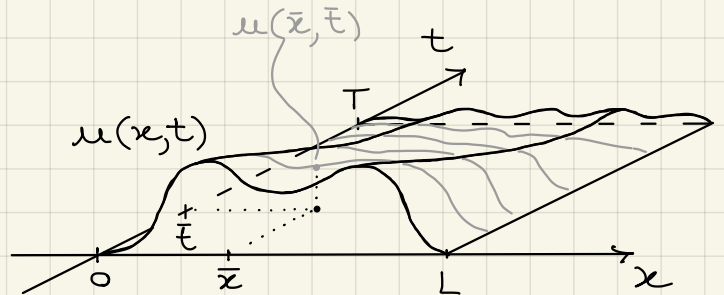
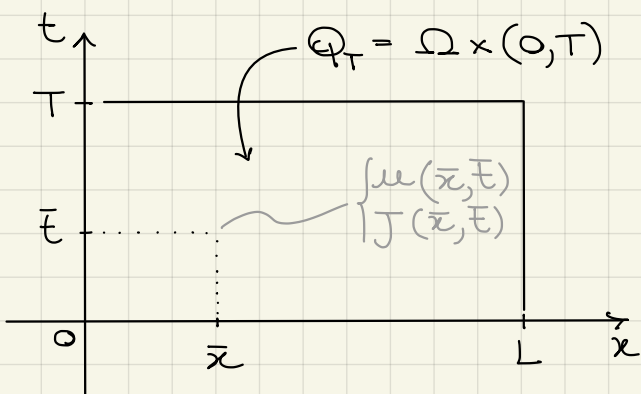
$J$  is very badly approximated (even if  $u$  is well approximated).  
Even using a higher order for the piecewise polynomial interpolation for  $u_n$  will improve  $J_n$  but won't make it continuous.

When approaching a numerical differential problem (with standard techniques) one has to choose whether he wants to sacrifice one variable or the other, depending on the application and final goal of the problem itself.

"DISPLACEMENT-BASED Formulation"  
( $u, u_n$ )

"MIXED / TWO-FIELD / HYBRID Methods"  
( $u_n, J_n, \{ \lambda_n, \mu_n \}$ )

advanced techniques that allow to retain both variables for a better physical description of the problem



We have to apply discretization to both  $x$  and  $t$ .

Note: discrete steps in time are typically addressed with  $\Delta t$ , while discrete steps in space are typically addressed with  $h$ .

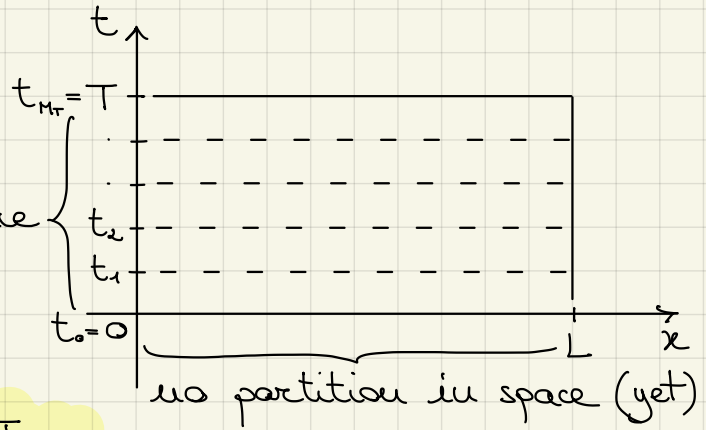


# Time semidiscretization

$$t_k = k\Delta t \quad \Delta t = \frac{T}{M_T} \quad M_T \gg 1$$

(uniform) partition in time

Information available only at discrete times.



$$\frac{u_{k+1} - u_k}{\Delta t} + \vartheta \frac{\partial J_{k+1}}{\partial x} + (1-\vartheta) \frac{\partial J_k}{\partial x} = \vartheta [g_{k+1} - K_{k+1} u_{k+1}] + (1-\vartheta) [g_k + K_k u_k]$$

$$k = 0, 1, \dots, M_T - 1 \quad x \in \Omega$$

$$\left\{ \begin{aligned} \left[ \frac{1}{\Delta t} + \vartheta K_{k+1} \right] u_{k+1} + \vartheta \left( \frac{\partial J_{k+1}}{\partial x} - g_{k+1} \right) &= \left[ \frac{1}{\Delta t} - (1-\vartheta) K_k \right] u_k - (1-\vartheta) \left( \frac{\partial J_k}{\partial x} + g_k \right) \\ J_{k+1} &= V_{k+1} + D \frac{\partial u_{k+1}}{\partial x} \end{aligned} \right.$$

$$t_k \rightarrow t_{k+1} \quad k = 0, 1, \dots, M_T - 1 \quad x \in \Omega$$

let's choose  $\vartheta = 1$  to simplify the notation:

$$\begin{aligned} (\tilde{P}) \quad & \left\{ \begin{aligned} \frac{\partial J_{k+1}}{\partial x} + \sigma_{k+1} u_{k+1} &= \frac{u_k}{\Delta t} + g_{k+1} & x \in (0, L) = \Omega \\ J_{k+1} &= V_{k+1} + D \frac{\partial u_{k+1}}{\partial x} & \sigma_{k+1} = \frac{1}{\Delta t} + K_{k+1} > 0 \\ \gamma_{k+1} \vec{J}_{k+1} \cdot \vec{n} &= \alpha_{k+1} u_{k+1} + \beta_{k+1} & x = 0; x = L \end{aligned} \right. \end{aligned}$$

$$u_0(x) \text{ given } \geq 0 \quad \forall x \in \bar{\Omega}$$

$\hookrightarrow f_1 \geq 0$  and what about future times?

## Monotone operators

$$V := C^2(\Omega) \cap C^1(\bar{\Omega})$$

two-times differentiable inside the domain

one-time differentiable at the boundary

$$\begin{aligned} Lu(x) &:= \frac{\partial J(u(x))}{\partial x} + \sigma(x) u(x) \\ &= \left[ \frac{\partial J(\cdot)}{\partial x} + \sigma(x) \cdot (\cdot) \right] u(x) \end{aligned}$$

for  $u \in V$

$$\Rightarrow Lu = f \iff u = L^{-1}f$$

"Placeholder" notation



In terms of linear algebra:  $\underline{L}u = \underline{f} \leftrightarrow u = \underline{L}^{-1} \underline{f}$

Def (Inverse-Monotonicity):

Let  $w \in V$  be such that

$$\left. \begin{array}{l} Lw(x) \geq 0 \quad \forall x \in \Omega \\ w(x) \geq 0 \quad \forall x \in \partial\Omega \end{array} \right\} \Rightarrow w(x) \geq 0 \quad \forall x \in \bar{\Omega}$$

Notation Clarification:

- $\vec{u}$  refers to a 3D vector (i.e.  $\vec{u} \in \mathbb{R}^3$ )
- $u$  refers to any nD vector (i.e.  $u \in \mathbb{R}^n \quad \forall n \geq 1$ )
- $U$  refers to any nxn matrix (i.e.  $U \in \mathbb{R}^{n \times n} \quad \forall n \geq 1$ )

Then  $L$  is said to be inverse-monotone.

Is this relevant for our applications? YES

The conditions to apply inverse monotonicity are almost always granted in physical applications of our interest ("u" is concentration, temperature etc.) so our variable will never take negative values.

Def (Maximum Principle):

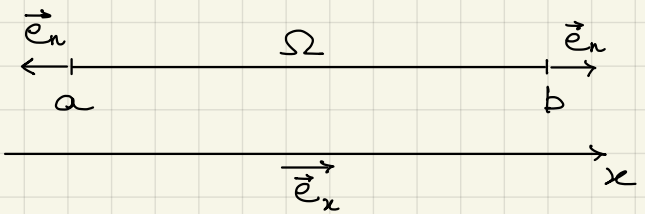
Let  $L$  be inverse monotone. Then  $L$  satisfies a maximum principle if

$$Lw(x) = 0 \quad x \in \Omega \Rightarrow \min_{x \in \partial\Omega} \{w(x), 0\} \leq w(x) \leq \max_{x \in \partial\Omega} \{w(x), 0\} \quad x \in \bar{\Omega}$$

$\parallel$   $W_{min}$   $\parallel$   $W_{max}$

This property ensures that our variable will always be bounded.

$$\begin{array}{l} Lu = f \quad \text{in } \Omega \\ B_{\partial\Omega} u = 0 \quad \text{on } \partial\Omega \end{array} \leftrightarrow \begin{cases} \frac{\partial J(u)}{\partial x} + \sigma u = f & \text{in } \Omega \\ J(u) = Vu - D \frac{\partial u}{\partial x} & \text{in } \Omega \\ \int_{\partial\Omega} \vec{J}(u) \cdot \vec{n} = \alpha_{\partial\Omega} u - \beta_{\partial\Omega} & \text{on } \partial\Omega \end{cases}$$



$$f(x) \geq 0 \quad \sigma(x) \geq 0 \quad x \in \bar{\Omega}$$

Assume now  $\begin{cases} \gamma_{\partial\Omega} = 0 \\ \alpha_{\partial\Omega} = 1 \\ \beta_{\partial\Omega} = 0 \end{cases} \Rightarrow \begin{cases} \frac{\partial J(u)}{\partial x} + \sigma u = f & \text{in } (a,b) \\ J(u) = Vu - D \frac{\partial u}{\partial x} & \text{"} \\ u = 0 & \text{at } x=a, x=b \end{cases} \quad (\tilde{P})$

Homogeneous Dirichlet Boundary conditions

We will now describe the numerical method to solve this boundary value problem  $(\tilde{P})$ .

## Finite Elements Method

Formal Steps to introduce the method:

1)  $\phi \cdot \left( \frac{\partial J}{\partial x} + \sigma u \right) = \phi f \quad \phi = \phi(x) \neq 0$

2)  $\int_a^b \phi \cdot \left( \frac{\partial J}{\partial x} + \sigma u \right) = \int_a^b \phi f \quad \forall \phi$

$\frac{\partial(\phi J)}{\partial x} = \frac{\partial \phi}{\partial x} \cdot J + \phi \cdot \frac{\partial J}{\partial x}$

3)  $\int_a^b \frac{\partial}{\partial x} (\phi J) - \int_a^b J \frac{\partial \phi}{\partial x} + \int_a^b \phi \sigma u = \int_a^b \phi f \quad \forall \phi \in C^1$

Assume  $\phi(a) = \phi(b) = 0 \sim u(a) = u(b) = 0$

Then  $\int_a^b \frac{\partial}{\partial x} (\phi J) = \phi(b)J(b) - \phi(a)J(a) = 0$

$\Rightarrow$  Weak formulation of the BVP  $(\tilde{P})$  is:

find  $u(x) \in X = \{w : (a,b) \rightarrow \mathbb{R} \mid \dots w(a) = w(b) = 0\}$  so that:

(w)  $\int_a^b \left( D \frac{\partial u}{\partial x} - Vu \right) \frac{\partial \phi}{\partial x} + \int_a^b \sigma u \phi = \int_a^b \phi f \quad \forall \phi \in X$

integral problem while  $(P)$  was a differential problem

$\phi$  is called "test function"

there might be some further properties

Short notation: (w)  $B(\phi, u) = F(\phi) \quad \forall \phi \in X$

bilinear form

linear form (linear functional)

Since  $\phi$  belongs to the same space as  $u$  we can arbitrarily set:

$$\phi = u$$

$$B(u, u) = F(u)$$

$$\int_a^b \left[ D \frac{\partial u}{\partial x} \frac{\partial u}{\partial x} - V u \frac{\partial u}{\partial x} + \sigma u u \right] = \int_a^b u f$$

$$\int_a^b \left[ D \left( \frac{\partial u}{\partial x} \right)^2 - V u \frac{\partial u}{\partial x} + \sigma u^2 \right] = \int_a^b u \cdot f$$

$$\int_a^b \left[ D \left( \frac{\partial u}{\partial x} \right)^2 - \frac{V}{2} \frac{\partial (u^2)}{\partial x} + \sigma u^2 \right] = \int_a^b u \cdot f$$

$$\int_a^b \left[ D \left( \frac{\partial u}{\partial x} \right)^2 - \frac{1}{2} \frac{\partial (V u^2)}{\partial x} + \frac{u^2}{2} \frac{\partial V}{\partial x} + \sigma u^2 \right] = \int_a^b u f, \quad V(x) \in C^1$$

$$\frac{1}{2} [V(b) u^2(b) - V(a) u^2(a)] = 0$$

$$\implies \int_a^b \left[ D \left( \frac{\partial u}{\partial x} \right)^2 + \left( \sigma + \frac{1}{2} \frac{\partial V}{\partial x} \right) u^2 \right] dx = \int_a^b u f dx$$

Assume that:

1.  $D(x) \geq D_{\min} > 0 \quad \forall x \in \bar{\Omega}$
2.  $\sigma(x) + \frac{1}{2} \frac{\partial V(x)}{\partial x} \geq \sigma_{\min} > 0 \quad \forall x \in \bar{\Omega}$

$$\implies B(u, u) \geq D_{\min} \int_a^b \left( \frac{\partial u(x)}{\partial x} \right)^2 dx + \sigma_{\min} \int_a^b u^2(x) dx$$

$$\text{but also } B(u, u) = F(u) = \int_a^b f(x) u(x) dx$$

Let us introduce the following space of functions

$$L^2(\Omega) := \left\{ w : \Omega \rightarrow \mathbb{R} \mid \int_{\Omega} w^2(x) dx < +\infty \right\}$$

which is an example of Hilbert space.

A property of Hilbert spaces is that functions can be treated as vectors: a scalar product between functions is therefore defined through the use of a norm.

For the  $L^2$  space this norm is:  $\|w\|_{L^2} = \sqrt{\int_{\Omega} w^2(x) dx} < +\infty$

So assuming now that  $u$  (and  $\phi$ ) and its derivative belong to  $L^2$  we can write:

$$B(u, u) \geq D_{\min} \left\| \frac{\partial u}{\partial x} \right\|_{L^2}^2 + \sigma_{\min} \|u\|_{L^2}^2$$

This assumption adds a specification to the regularity of  $u$  and therefore to the space of functions  $X$ :

$$X = \left\{ w: \Omega \rightarrow \mathbb{R} \mid \begin{array}{l} \text{a) } w \in L^2(\Omega) \\ \text{b) } \frac{\partial w}{\partial x} \in L^2(\Omega) \\ \text{c) } w(x) \Big|_{\partial\Omega} = 0 \end{array} \right\} \#$$

A function space with such properties is typically indicated with the symbol:

$$X = H_0^1(\Omega)$$

and its norm is defined as:  $\|w\|_{H_0^1} = \sqrt{\|w\|_{L^2}^2 + \left\| \frac{\partial w}{\partial x} \right\|_{L^2}^2} < +\infty$

$$\rightarrow B(u, u) \geq \underbrace{\min\{D_{\min}, \sigma_{\min}\}}_{\alpha_0} \cdot \|u\|_{H_0^1}^2$$

$$\alpha_0 \|u\|_{H_0^1}^2 \leq \int_a^b f u \, dx$$

If we also assume  $f \in L^2(\Omega)$  then we can apply the Cauchy-Schwartz inequality:

$$\|f \cdot u\|_{L^1} = \int_a^b f u \, dx \leq \|f\|_{L^2} \|u\|_{L^2}$$

Also by definition:  $\|u\|_{H_0^1} \geq \|u\|_{L^2}$ .

We can finally write:

$$\alpha_0 \|u\|_{H_0^1}^2 \leq \int_a^b f u \, dx \leq \|f\|_{L^2} \|u\|_{L^2} \leq \|f\|_{L^2} \|u\|_{H_0^1}$$

$$\Rightarrow \boxed{\|u\|_{H_0^1} \leq \frac{\|f\|_{L^2}}{\alpha_0}} \rightarrow \begin{array}{l} \text{a priori estimate of} \\ \text{the solution} \\ \downarrow \\ \text{stability estimate} \\ \downarrow \\ \text{well-posedness} \\ \downarrow \\ \text{continuous on the data!} \end{array}$$

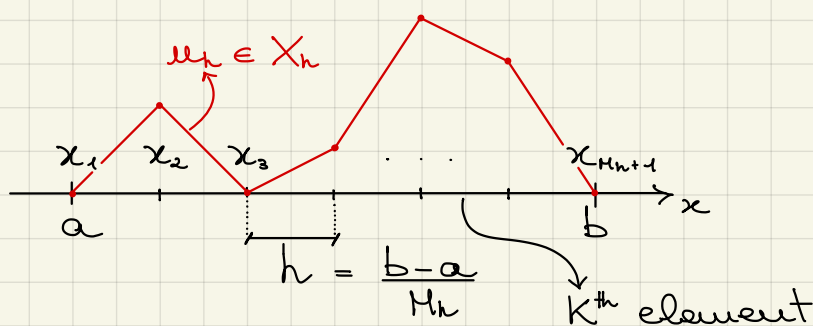
This estimate was obtained through several assumptions. Let's summarize them all together:

- $u \in H_0^1(\Omega)$
- $V \in C^1(\Omega)$  ( $V(x)$  is differentiable and  $\frac{\partial V}{\partial x} \in L^\infty(\Omega)$ )
- $D(x) \geq D_{\min} > 0 \quad \forall x \in \bar{\Omega}$
- $\mathcal{O}(x) + \frac{1}{2} \frac{\partial V(x)}{\partial x} \geq \mathcal{O}_{\min} > 0 \quad \forall x \in \bar{\Omega}$
- $f(x) \geq 0$  and  $f \in L^2(\Omega)$

Only under these assumptions we can say that the weak formulation of problem (P) admits a unique solution that depends with continuity on the data according to the above stability estimate.

(P) strong formulation  $\longrightarrow$  (W) weak formulation  
 2nd order BVP " find  $u \in X = H_0^1(\Omega)$  so that  
 $F(x, d) = 0 \Leftrightarrow B(u, \phi) - F(\phi) = 0 \quad \forall \phi \in X$ "  
 $H_0^1$  is a Sobolev space

We need to apply a discretization of  $\Omega$  to the analytical problem (W) in order to obtain a computable problem. In this way we will define the finite element method.



$$F(u, d) \longrightarrow F_h(u_n, d_h)$$

$$X \longrightarrow X_n$$

$$\dim(X_n) = M_n - 1$$

$$u_n(a) = u_n(b) = 0$$

$$u_n \in C^0(\bar{\Omega})$$

$$\forall K \in \mathcal{T}_n: u_n|_K \in P_1(K)$$

union of all  $K$  elements  $\downarrow$  polynomials of degree 1 over  $K$   
 "restriction operator"

$M_n$ : # of elements

$h$ : discretization parameter

An element is a subinterval of the entire domain.

$$M_n \longrightarrow +\infty \quad h \longrightarrow 0$$

The space of functions  $X_h$  so defined is called "finite element space of degree  $l$  associated with the partition  $T_h$ ".

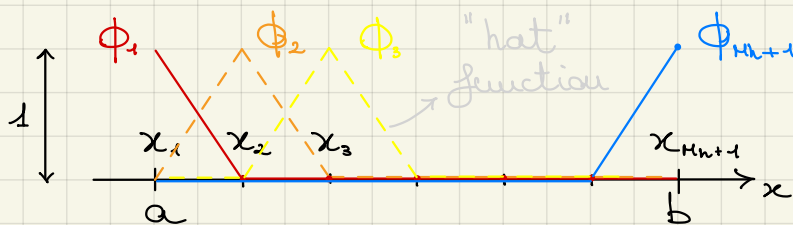
The only needed information to derive the entire  $u_h$  is the value at each node.

The difference between finite elements and finite differences methodologies is that the former returns the union of polynomial functions as approximated solution, while the latter only returns a set of values that approximate the solution just at the nodes.

$$X_h = \text{span}\{\phi_j\}_{j=1}^{N_h} \quad N_h = \dim(X_h) = M_h - 1$$

$$\rightarrow \forall u_h \in X_h: u_h(x) = \sum_{j=1}^{N_h} u_j \phi_j(x)$$

Let's see how the basis functions  $\phi_j(x)$  are made:



Since  $\phi_1$  and  $\phi_{M_h+1}$  are set by boundary conditions ( $u_1 = u(a)$  and  $u_{M_h+1} = u(b)$ ) they are not accounted for in the aforesaid expression (so  $\phi_2$  is actually  $\phi_1$  in the previous notation).

A more general notation would be:

$$X_h = \text{span}\{\phi_j\}_{j=1}^{N_h+2} \rightarrow u_h(x) = \sum_{j=1}^{N_h+2} u_j \phi_j(x)$$

where for  $j=1$  and  $j=N_h+2$  the polynomial is given by the boundary conditions.

We can finally define the finite element method for our weak formulation of the BVP:

find  $u_h \in X_h$  so that  $B(u_h, \phi_i) = F(\phi_i) \quad i = 1, \dots, N_h+2$

$$u_h(x) = \sum_{j=1}^{N_h+2} u_j \phi_j(x)$$

Finite Element equations:  $\sum_{j=1}^{N_h+2} u_j B(\phi_j, \phi_i) = F(\phi_i)$



In terms of linear algebra:

$$\underline{B} \underline{u} = \underline{F}$$

i row  
j column

where  $(\underline{B})_{ij} = B(\phi_j, \phi_i)$

$$(\underline{F})_i = F(\phi_i)$$

→ The finite element method entails the resolution of  $N_h + 2$  equations (plus the computation of the stiffness matrix  $\underline{B}$  and the load vector  $\underline{F}$ ) to obtain the vector of nodal unknowns  $\underline{u}$

This linear algebraic system is solvable only if  $\underline{B}$  is invertible (non-singular), which is granted by the assumptions made in our preliminary discussion.

In other words, because of these  $\square$  assumptions, it can be demonstrated that  $\underline{B}$  is positive definite.

The one discussed so far is a FEM of degree 1, since the polynomials used for interpolation were of first order (linear).

A higher degree FEM can also be used with higher order polynomials, possibly increasing the accuracy of the result.

### Error estimate

find  $u \in X$  so that

$$(W) \quad B(u, \phi) = F(\phi) \quad \forall \phi \in X$$

→ non-computable

find  $u_h \in X_h \subset X$  so that

$$(W_h) \quad B(u_h, \phi_h) = F(\phi_h) \quad \forall \phi_h \in X_h$$

→ computable

BUT we know that

$$\|u\|_X \leq \frac{\|f\|_{L^2}}{\alpha_0}$$

$$e_h = u - u_h$$

→ computable or not?

Theorem. Assume that  $u \in H^2(\Omega) \cap H_0^1(\Omega)$ .

$$\text{Then } \|u - u_h\|_{H_0^1}^{\uparrow} \leq c h^{\uparrow} \|u\|_{H^2}$$

Theorem. Assume that  $u \in H^2(\Omega) \cap H_0^1(\Omega)$ .

$$\text{Then } \|u - u_h\|_{L^2}^{\downarrow} \leq c h^{2\downarrow} \|u\|_{H^2}$$

The two theorems are consistent with each other, since the norm in  $H^1$  is always greater, by definition, than the norm in  $L^2$ . Therefore the error estimate is less restrictive (1st order convergence) for the  $H^1$  norm while it is more restrictive (2nd order convergence) for the  $L^2$  norm.

The weak formulation we wrote for our problem ( $\tilde{P}$ ):

$$(W) \int_a^b \left( D \frac{\partial u}{\partial x} - Vu \right) \frac{\partial \phi}{\partial x} + \int_a^b \sigma u \phi = \int_a^b \phi f$$

is called displacement-based weak formulation, as it uses "u" as a variable and loses the dependency on "J".

This, as we have already highlighted, will cause some issues when retrieving the values for J.

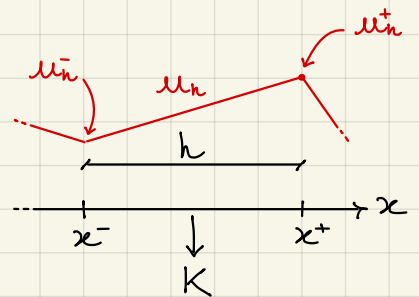
$$J(u) = Vu - D \frac{\partial u}{\partial x}$$

1st order FEM

$$J_h = V \tilde{u}_h - D \frac{\partial \tilde{u}_h}{\partial x}$$

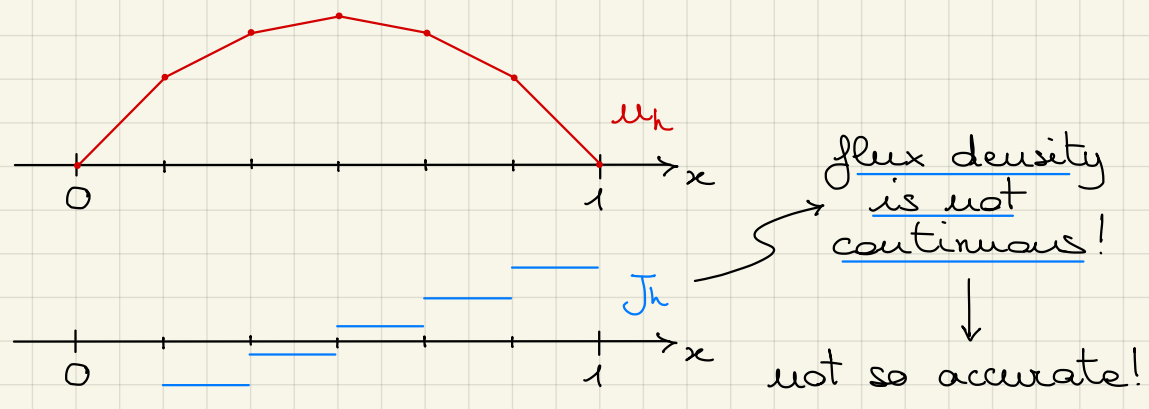
where

$$\begin{cases} \frac{\partial \tilde{u}_h}{\partial x} \Big|_K = \frac{u_n^+ - u_n^-}{h} \\ \tilde{u}_h \Big|_K = \frac{\int_K u_h dx}{h} = \frac{u_n^- + u_n^+}{2} \end{cases}$$



$$\Rightarrow J_h \Big|_K = V \frac{u_n^+ + u_n^-}{2} - D \frac{u_n^+ - u_n^-}{h} \quad \forall K \in \mathcal{T}_h$$

If for example  $V(x) = D(x) = \sigma(x) = f(x) \equiv 1$ ,  
 $(a, b) = (0, 1)$  and  $N_h = 5$ :





Another issue associated with the FEM we have analyzed resides in the "a priori" estimate of the solution:

$$\|u\|_{H^1} \leq \frac{\|f\|_{L^2}}{\alpha_0}$$

where  $\alpha_0 = \min\{D_{\min}, \sigma_{\min}\}$ ,  $D_{\min} \leq D(x)$ ,  $\sigma_{\min} \leq \sigma(x) + \frac{1}{2} \frac{\partial V(x)}{\partial x}$

This estimate also holds for the numerical solution:

$$\|u_h\|_{H^1} \leq \frac{\|f\|_{L^2}}{\alpha_0}$$

It is evident that, as  $\alpha_0$  decreases, the estimate allows  $u_h$  to reach very high values.

It then happens that for very low  $\alpha_0$  the numerical solution suffers from spurious oscillations unless a very narrow discretization parameter  $h$  is adopted. These oscillations of course make the computed solution unreliable from a physical standpoint.

Let us now introduce the Péclet Numbers:

$$Pe_{ad} = \frac{|V|h}{2D}$$

Péclet number associated with advection

$$Pe_{\text{reac}} = \frac{\sigma h^2}{6D}$$

Péclet number associated with reaction

It can be demonstrated that if either of the two Péclet numbers is greater than one then the numerical solution is affected by spurious oscillations.

There exist "artificial" ways to decrease the Péclet number (whichever is too high) without reducing the discretization parameter, which typically requires too high computational costs.

Since both numbers depend on  $D$ , these techniques generally consist of a (moderate) increase of the diffusion parameter:

$$D \longrightarrow D_{\text{eff}} = D + \text{"extra diffusivity"} \\ = D (1 + G(Pe))$$

E.g.:  $G(\text{Pe}_{ad}) = \text{Pe}_{ad} \longrightarrow D_h = D(1 + \text{Pe}_{ad}) =$

"Upwind stabilization"

$$\tilde{\text{Pe}}_{ad} = \frac{|V|h}{2D_h} = \frac{|V|h}{2D(1 + \text{Pe}_{ad})}$$

$$= \frac{\text{Pe}_{ad}}{1 + \text{Pe}_{ad}} < 1 \text{ always!}$$

The result will be perturbed but it won't have oscillations!

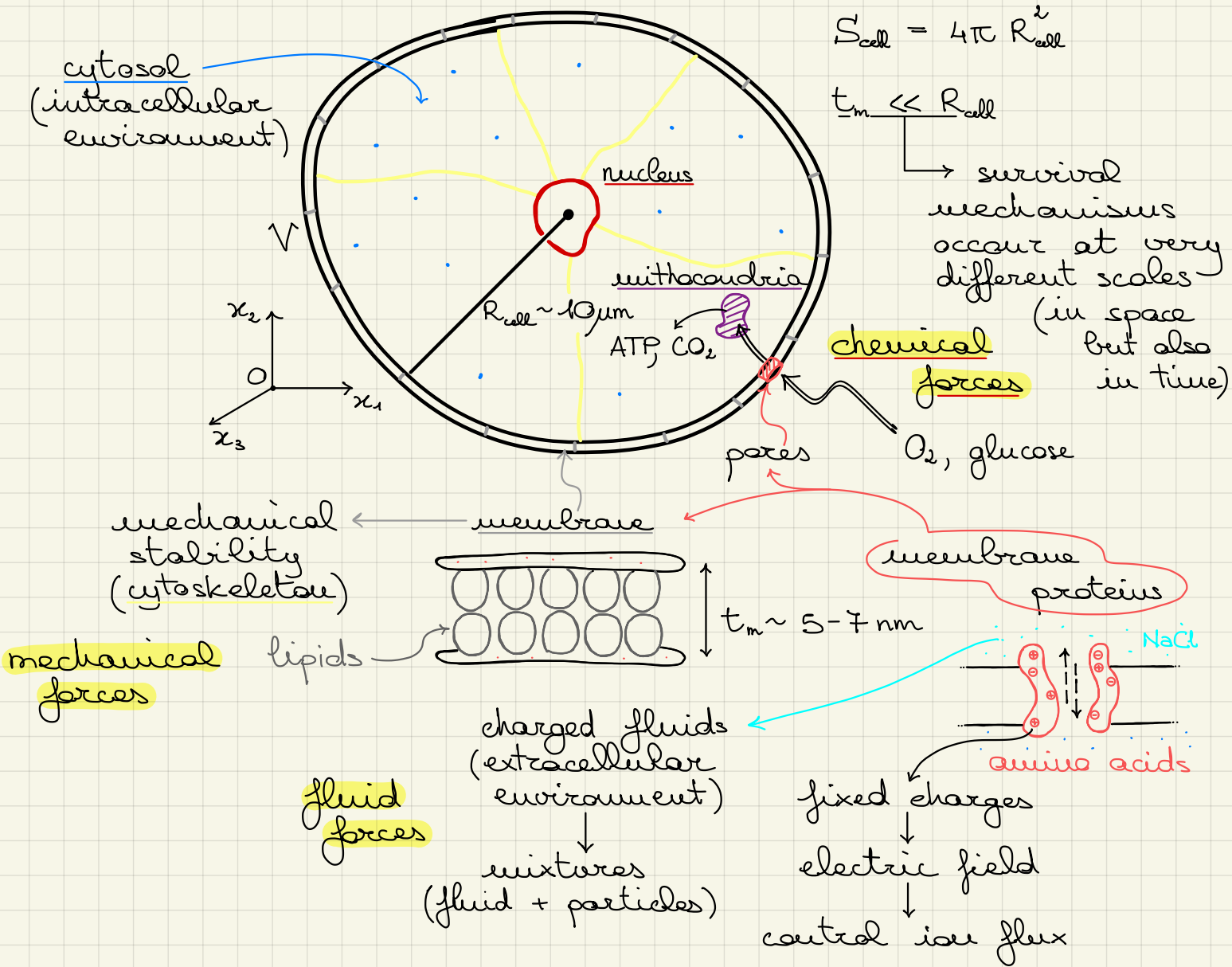
Note: the condition on the Péclet numbers is reminiscent of the absolute stability condition of the forward Euler method.

$$\text{Pe}_{ad} = \frac{|V|h}{2D} < 1 \longrightarrow h < \frac{2D}{|V|}$$

$$\Delta t < \frac{2}{\lambda} \quad \left. \begin{array}{l} \updownarrow \\ \lambda \approx \frac{|V|}{D} \end{array} \right\}$$

# Cellular Biology and Electrophysiology

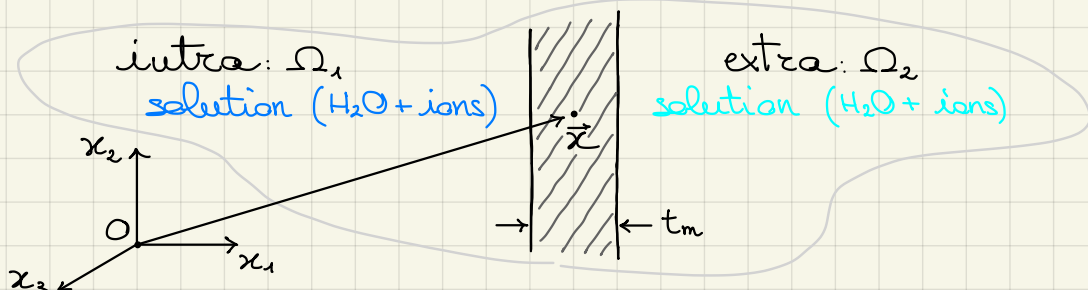
## The Cell



## HOMEOSTASIS

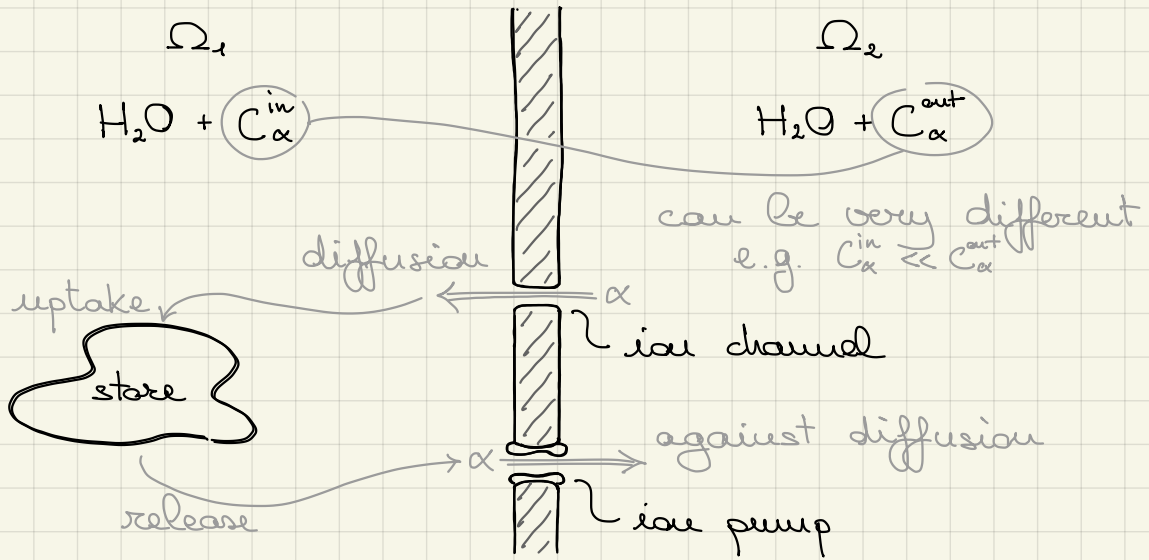
balance laws + constitutive equations in a cell

We will focus on the electrochemical activity of the cell and forget about other mechanisms to simplify the study of such complex system.

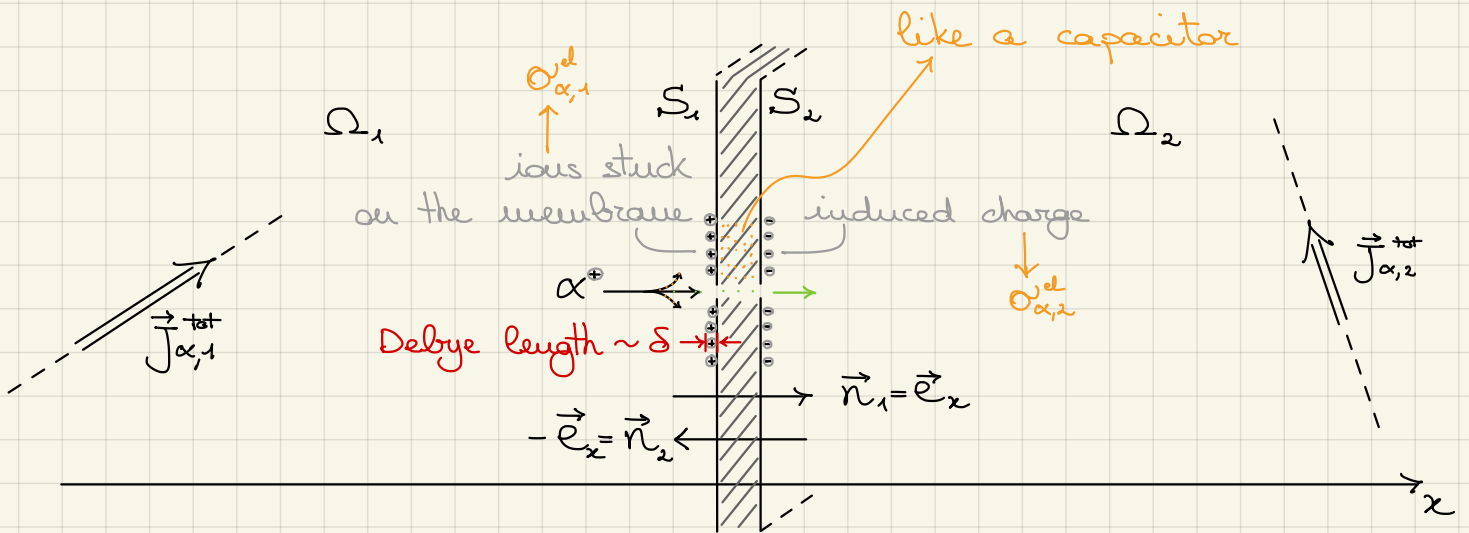


Molar density of  $\alpha$  ion species:  $C_\alpha = \left[ \frac{\text{mol}}{\text{m}^3} \right] = [\text{mM}]$

$\frac{\text{mol}}{\text{m}^3} = \frac{\text{mol}}{10^3 \text{dm}^3} = 10^{-3} \frac{\text{mol}}{\text{l}} = 10^{-3} \text{M}$  where 1 mole = 1M "molar" liter



e.g.:  $\alpha$  is  $\text{Ca}^{2+}$  and store is endoplasmic reticulum  
 $\hookrightarrow$  noxious at high concentrations



$$\vec{j}_\alpha^{\text{tot}} \cdot \vec{n} = j_\alpha^{\text{tot}} = \left[ \frac{\text{A}}{\text{m}^2} \right] = \left[ \frac{\text{C}}{\text{s m}^2} \right]$$

$$j_\alpha^{\text{tot}} = j_\alpha^{\text{capacitive}} + j_\alpha^{\text{conductive}}$$

$\frac{\partial \rho_\alpha^{\text{dl}}}{\partial t}$  variation of charge density

$$\rho_{\alpha,1}^{\text{dl}}(t) + \rho_{\alpha,2}^{\text{dl}}(t) = 0 \quad \forall t$$

Continuity equation:

$$\text{rot}(\vec{H}) = \vec{j} + \frac{\partial \vec{D}}{\partial t} = \vec{j}_{\text{tot}} \rightarrow \vec{\nabla} \cdot \vec{j} = 0$$

conduction
displacement

$$\text{KCL: } -j_{\alpha,1}^{\text{tot}} + \frac{\partial \rho_{\alpha,1}^{\text{el}}}{\partial t} + j_{\alpha,1}^{\text{cond}} = 0$$

$$\text{KCL: } -j_{\alpha,2}^{\text{tot}} + \frac{\partial \rho_{\alpha,2}^{\text{el}}}{\partial t} + j_{\alpha,2}^{\text{cond}} = 0$$

$$j_{\alpha,1}^{\text{tot}} = \vec{j}_{\alpha,1}^{\text{tot}} \cdot \vec{n}_1$$

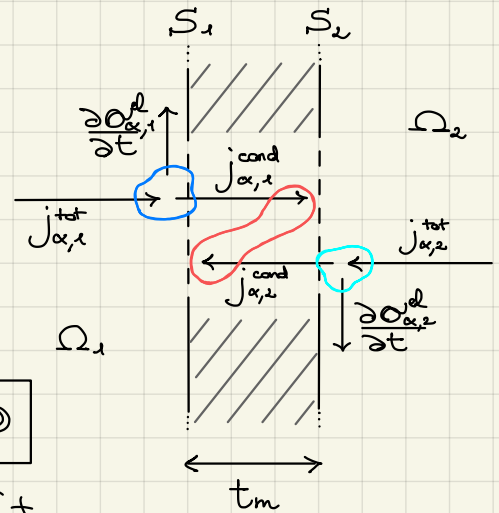
$$j_{\alpha,2}^{\text{tot}} = \vec{j}_{\alpha,2}^{\text{tot}} \cdot \vec{n}_2$$

$$\vec{\nabla} \cdot \vec{j}_{\alpha}^{\text{tot}} = 0$$

continuity equation

$$\rho_{\alpha,1}^{\text{el}} + \rho_{\alpha,2}^{\text{el}} = 0$$

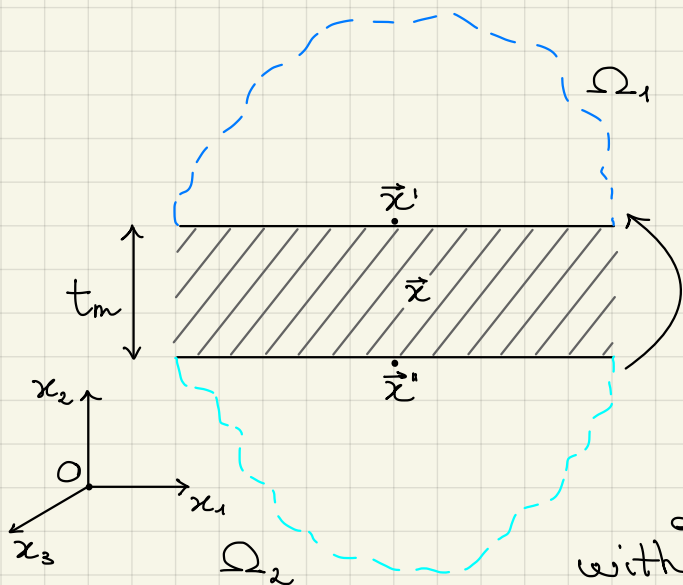
electroneutrality of the membrane



$$\Rightarrow -j_{\alpha,1}^{\text{tot}} - j_{\alpha,2}^{\text{tot}} + \underbrace{j_{\alpha,1}^{\text{cond}} + j_{\alpha,2}^{\text{cond}}}_{= 0 \text{ because of KCL}} = 0 \quad \forall t$$

$$\Rightarrow j_{\alpha,1}^{\text{tot}} = -j_{\alpha,2}^{\text{tot}}$$

(no absorption of ions within the membrane)



In the cell timescale it is typically:

$$\text{rot}(\vec{E}) = \frac{\partial \vec{B}}{\partial t} \approx 0 \rightarrow \vec{E} = -\vec{\nabla} \psi \quad \text{quasi-static approx.}$$

$$V_m(t) := \psi^{\text{in}}(t) - \psi^{\text{out}}(t) \quad [\text{Volt}]$$

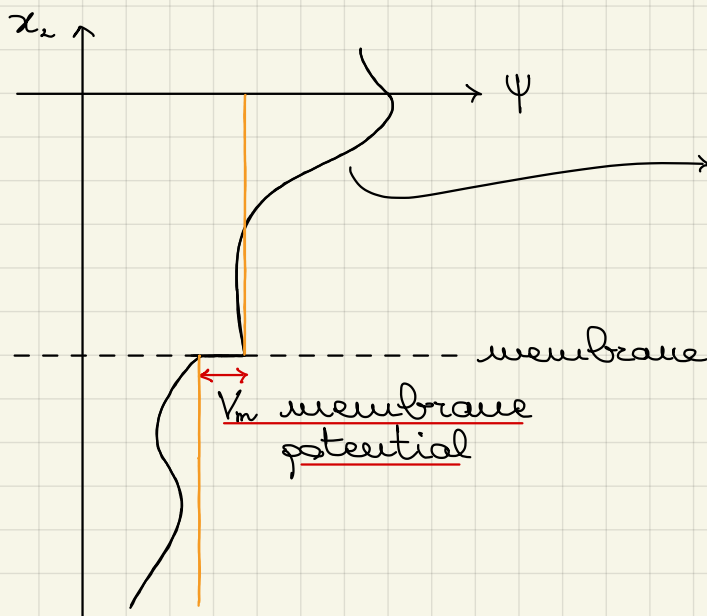
The electric field can be considered irrotational. In other words, we are dealing with stationary fields.

$$\psi^{\text{in}}(t) := \psi(\vec{x}', t) \quad \psi^{\text{out}}(t) := \psi(\vec{x}'', t)$$

But  $\vec{x}' \approx \vec{x} \approx \vec{x}''$  because of the short thickness of the membrane with respect to the surroundings.

However  $\psi^{\text{in}} \neq \psi^{\text{out}}$  certainly.

If we then let  $t_m \rightarrow 0$  we obtain a discontinuous function of  $\psi$  where the discontinuity corresponds to the membrane and is equal in amplitude to  $V_m$ .



the electric potential indeed changes within and outside the cell, however from now on we will assume it to be constant

$$\psi(\vec{x}, t) = \psi_x(t)$$

$$V_m(\vec{x}, t) = \psi_x^{\text{in}}(t) - \psi_x^{\text{out}}(t)$$

local membrane potential

Definition:  $j_\alpha^{\text{cap}}(\vec{x}, t) := C_m \frac{\partial V_m(\vec{x}, t)}{\partial t}$

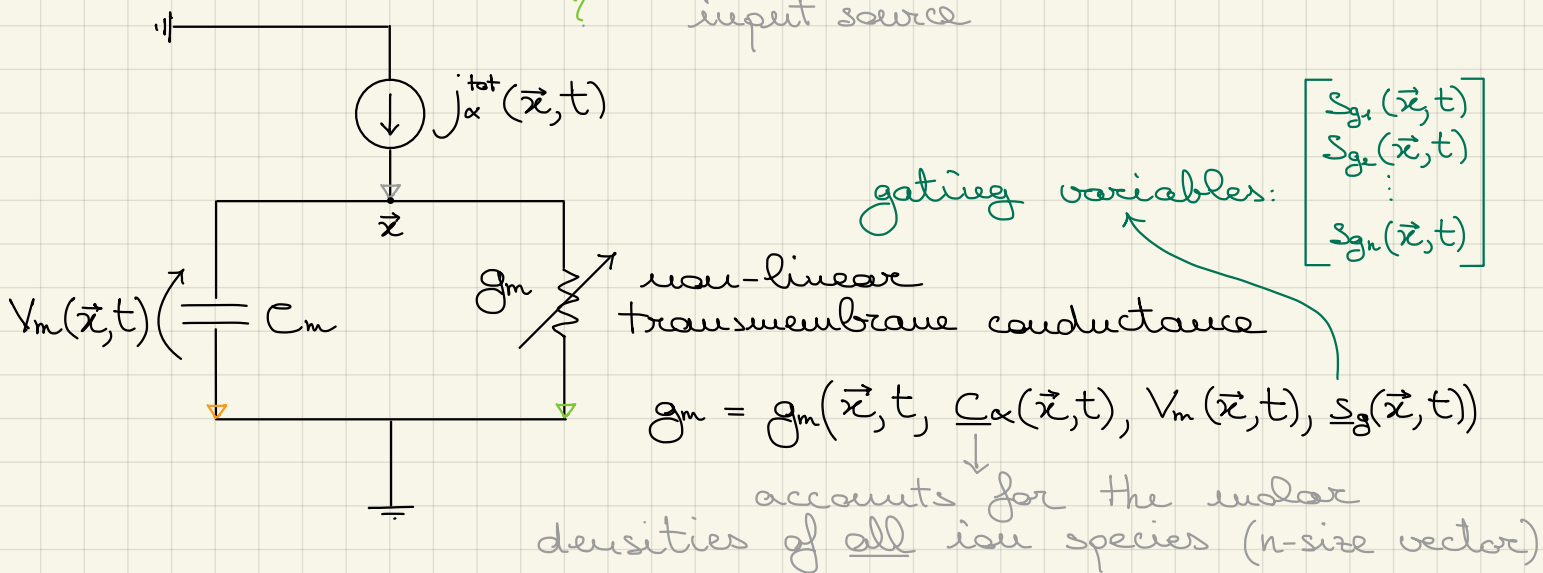
↓  
specific capacitance

$$C_m = \frac{\epsilon_m}{t_m} = \left[ \frac{\text{F}}{\text{m}^2} \right]$$

dielectric constant of the membrane

From experiments (Hodgkin-Huxley):  $\frac{\epsilon_m}{t_m} = C_m \sim 10^{-2} \frac{\text{F}}{\text{m}^2}$

$$C_m \frac{\partial V_m(\vec{x}, t)}{\partial t} + j_\alpha^{\text{cond}}(\vec{x}, t) = j_\alpha^{\text{tot}}(\vec{x}, t) \quad \forall \vec{x} \in \partial V \quad \forall t \in I_T$$



The gating variables  $s_{g_j} \in [0, 1]$  are probabilistic quantities that describe the likelihood of the membrane pores of being open or closed.

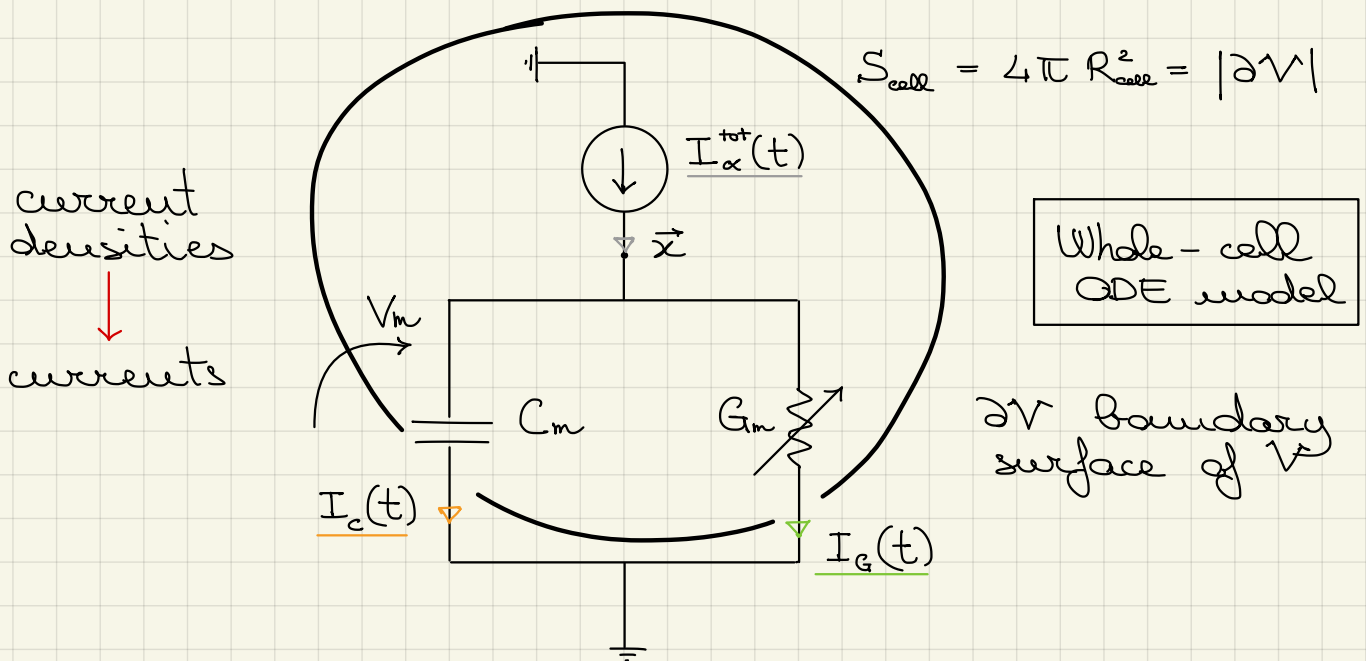


This is the equivalent mathematical model of a piece of membrane (local ODE model)

Note how for each fixed  $\vec{x}$  it resembles a Cauchy problem.

## Characterization of the transmembrane conductance

Let's apply all previous considerations, not just to a portion of the membrane ( $\vec{x}$  inside the membrane), but to the whole cell volume ( $\vec{x}$  at the center of the cell).



$$\vec{j}_\alpha^{\text{tot}}(\vec{x}, t) \cdot \vec{n} = j_\alpha^{\text{tot}}(\vec{x}, t) \longrightarrow \int_{S_{\text{cell}}} \vec{j}_\alpha^{\text{tot}}(\vec{x}, t) \cdot \vec{n} \, d\Sigma = I_\alpha^{\text{tot}}(t)$$

$$-I_\alpha^{\text{tot}}(t) + I_c(t) + I_g(t) = 0$$

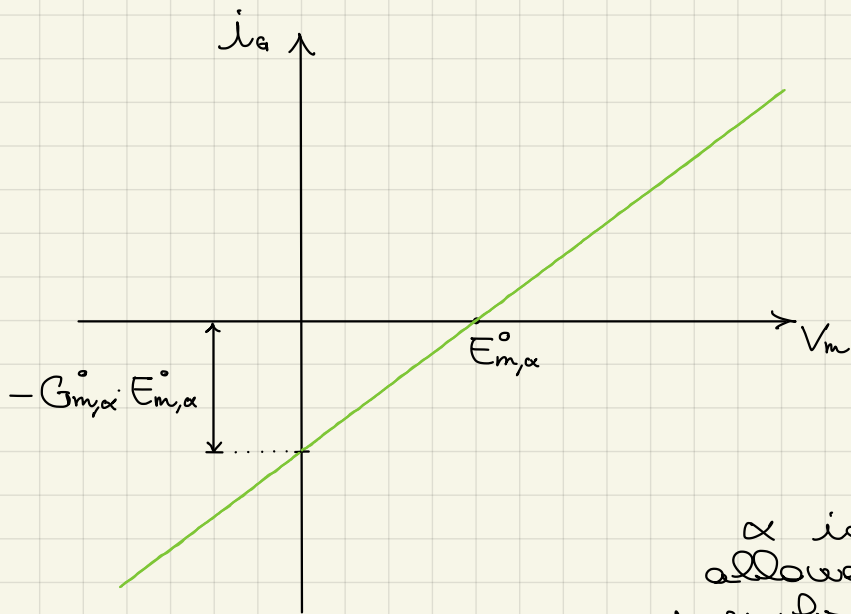
where  $I_c(t) = C_m \cdot \frac{dV_m(t)}{dt}$  with  $C_m = C_m \cdot S_{\text{cell}}$

$$I_g(t) = ?$$

linear resistor

e.g.:  $I_g(t) = G_{m,\alpha}(t, V_m(t)) [V_m(t) - E_{m,\alpha}(t, V_m(t))]$

1) linear resistor model:  $I_g(t) = G_{m,\alpha}^\circ [V_m(t) - E_{m,\alpha}^\circ]$



At steady-state:  $i_c = 0$

If  $V_m = E_{m,\alpha}^0$ :  $i_\alpha = 0$

$$I_\alpha^{\text{tot}} = 0$$

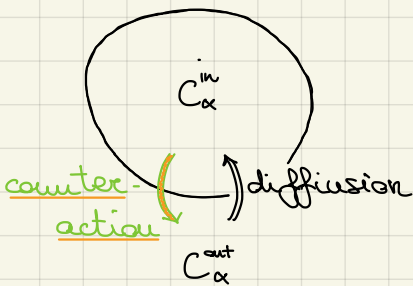
"thermodynamic equilibrium"

$\alpha$  ions are not allowed to cross the membrane, even if  $C_\alpha^{\text{in}} \neq C_\alpha^{\text{out}}$

This is achievable only for a special value of the membrane potential:

$$V_m = E_{m,\alpha}^0$$

which has to be associated to some form of work within the cell that opposes to the natural gradient diffusion of the ions.



$\Rightarrow$  Nernst potential associated with ion species  $\alpha$ :

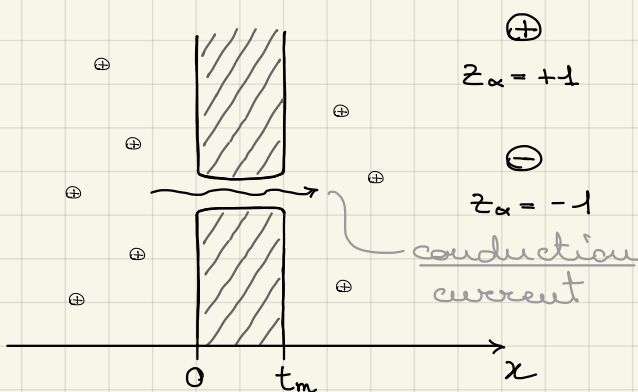
$$E_{m,\alpha}^0 := \frac{V_{th}}{z_\alpha} \ln\left(\frac{C_\alpha^{\text{out}}}{C_\alpha^{\text{in}}}\right)$$

where  $V_{th} := \frac{k_B T}{q}$  thermal voltage ( $V_{th} \approx 26 \text{ mV}$  @  $T = 300 \text{ K}$ )

$$k_B = 1,38 \cdot 10^{-23} \text{ J/K} \quad q = 1,602 \cdot 10^{-19} \text{ C}$$

$z_\alpha$ : chemical valence of ion species  $\alpha$

Demonstration of Nernst potential formula:



$\oplus$   
 $z_\alpha = +1$   
 $\ominus$   
 $z_\alpha = -1$

$$\int_\Omega \rho_\alpha = q z_\alpha N_{AV} C_\alpha = F z_\alpha C_\alpha$$

$$\int_\Omega q z_\alpha N_{AV} C_\alpha d\Omega = Q_\alpha$$

volumetric charge density

$N_{AV}$ : Avogadro number

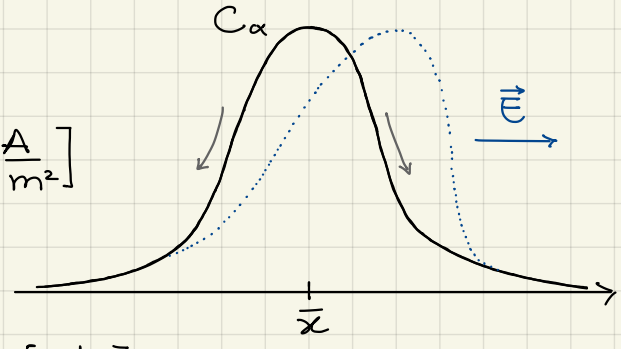
$F = q N_{AV}$ : Faraday constant



Ions can move by:

+ diffusion  $\vec{J}_{diff} = -F z_\alpha D_\alpha \vec{\nabla} C_\alpha$   
 current density  $= \left[ \frac{C}{mol} \cdot \frac{m^2}{s} \cdot \frac{1}{m} \cdot \frac{mol}{m^3} \right] = \left[ \frac{A}{m^2} \right]$

+ electric force  $\vec{F} = F z_\alpha C_\alpha \vec{E}$   
 force density  $= \left[ \frac{C}{mol} \cdot \frac{mol}{m^3} \cdot \frac{V}{m} \right] = \left[ \frac{N}{m^3} \right]$



electrical mobility  $\mu_\alpha^{el} = \left[ \frac{m^2}{V \cdot s} \right]$

$\vec{J}_{drift} = \mu_\alpha^{el} F |z_\alpha| C_\alpha \vec{E} = \sigma_\alpha \vec{E}$

where  $\sigma_\alpha = F |z_\alpha| C_\alpha \mu_\alpha^{el} = \left[ \frac{\Omega^{-1}}{m} \right]$   
 electrical conductivity

⇒ Nernst-Planck transport model:

$$\vec{J}_\alpha = F |z_\alpha| \mu_\alpha^{el} C_\alpha \vec{E} - F z_\alpha D_\alpha \vec{\nabla} C_\alpha$$

It is a model for ion electrodiffusion (drift-diffusion, like electrons and holes in a semiconductor)

Impose now  $J_\alpha = F |z_\alpha| \mu_\alpha^{el} C_\alpha E - F z_\alpha D_\alpha \frac{\partial C_\alpha}{\partial x} = 0$  (equilibrium)  
 $-\frac{\partial \psi}{\partial x} = E$

$|z_\alpha| \mu_\alpha^{el} C_\alpha \frac{\partial \psi}{\partial x} + z_\alpha D_\alpha \frac{\partial C_\alpha}{\partial x} = 0$

Einstein-Stokes-Smoluchowski equation:  $D_\alpha = \frac{V_{th} \mu_\alpha^{el}}{|z_\alpha|}$

$\frac{\partial C_\alpha}{\partial x} = - \frac{|z_\alpha|^2 C_\alpha}{z_\alpha V_{th}} \frac{\partial \psi}{\partial x}$

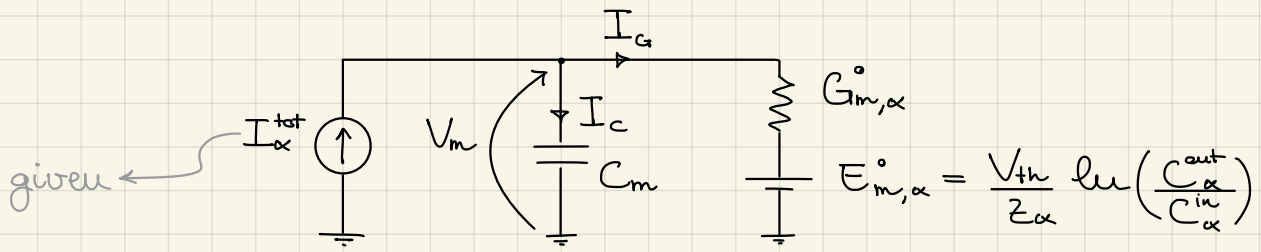
$\frac{\partial \psi}{\partial x} = - \frac{V_{th}}{z_\alpha} \frac{\partial C_\alpha}{\partial x} \frac{1}{C_\alpha} = - \frac{V_{th}}{z_\alpha} \frac{\partial \ln(C_\alpha)}{\partial x}$

$\psi = \psi(x)$  and  $C_\alpha = C_\alpha(x)$  since we are at equilibrium (steady-state, no time dependence).

$\psi(x) - \psi(0) = - \frac{V_{th}}{z_\alpha} \ln \left( \frac{C_\alpha(x)}{C_\alpha(0)} \right) =$  "chemical potential"

$x = t_m \implies V_m = \frac{V_{th}}{z_\alpha} \ln \left( \frac{C_\alpha^{out}}{C_\alpha^{in}} \right)$

# Whole-cell (linear resistor) model



$$I_c = C_m \frac{dV_m}{dt} \quad I_G = G_{m,\alpha} (V_m - E_{m,\alpha}^{\circ})$$

$$-I_{\alpha}(t) + C_m \frac{dV_m(t)}{dt} + G_{m,\alpha} (V_m(t) - E_{m,\alpha}^{\circ}) = 0$$

$$\left( R_{m,\alpha}^{\circ} = \frac{1}{G_{m,\alpha}^{\circ}} \right) \quad \frac{dV_m(t)}{dt} = \frac{I_{\alpha}^{\text{tot}}}{C_m} + \frac{E_{m,\alpha}^{\circ}}{R_{m,\alpha}^{\circ} C_m} - \frac{V_m(t)}{R_{m,\alpha}^{\circ} C_m} \rightarrow \tau$$

$$\begin{cases} \frac{dV_m}{dt} = -\frac{V_m}{\tau} + \frac{1}{\tau} [E_{m,\alpha}^{\circ} + I_{\alpha}^{\text{tot}} R_{m,\alpha}^{\circ}] & t \in I_T = (t_0, t_0 + T) \\ V_m(t_0) = V_m^{\circ} = E_{m,\alpha}^{\circ} \rightarrow \text{equilibrium at starting condition} \end{cases}$$

Let's now analyze the Nernst-Planck equation by modifying its expression in a more familiar form:

$$\vec{J}_{\alpha} = q \mu_{\alpha}^{\text{el}} |z_{\alpha}| n_{\alpha} \vec{E} - q z_{\alpha} D_{\alpha} \vec{\nabla} \cdot n_{\alpha} \quad n_{\alpha}(\vec{x}) = C_{\alpha}(\vec{x}) \cdot N_{\alpha}$$

$$1D \rightarrow \vec{e}_x \left[ q \mu_{\alpha}^{\text{el}} |z_{\alpha}| n_{\alpha}(x) E - q z_{\alpha} D_{\alpha} \frac{\partial n_{\alpha}(x)}{\partial x} \right]$$

$$J_{\alpha} = q z_{\alpha} \left[ \mu_{\alpha}^{\text{el}} \frac{|z_{\alpha}|}{z_{\alpha}} n_{\alpha} E - D_{\alpha} \frac{\partial n_{\alpha}}{\partial x} \right]$$

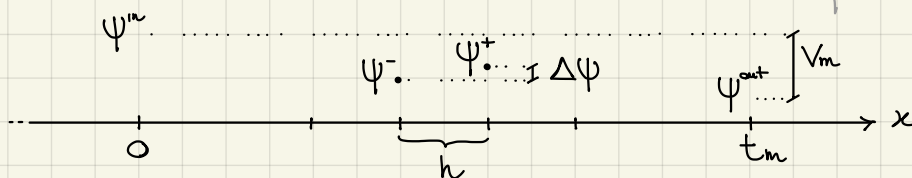
$$\rightarrow f_{\alpha} = \frac{J_{\alpha}}{q z_{\alpha}} = \left[ \frac{C}{s \cdot m^2 \cdot C} \right] = [m^2 \cdot s^{-1}]$$

$$f_{\alpha}(n_{\alpha}) = \underbrace{\mu_{\alpha}^{\text{el}} \frac{|z_{\alpha}|}{z_{\alpha}} E}_{\psi} n_{\alpha} - D_{\alpha} \frac{\partial n_{\alpha}}{\partial x} \quad \longleftrightarrow \quad J(u) = Vu - D \frac{\partial u}{\partial x}$$

Compute the Péclet number associated to this advection-diffusion problem:

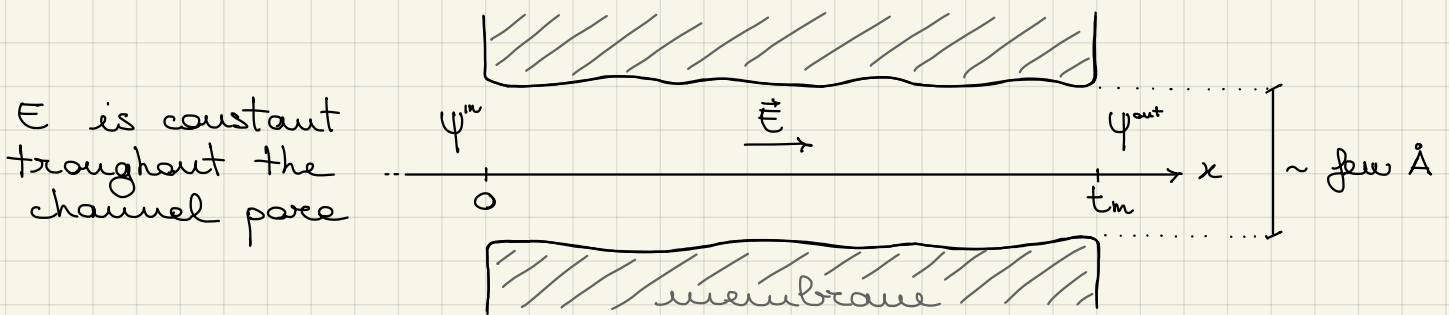
$$Pe_{\text{ad}} = \frac{h|V|}{2D} = \frac{h \mu_{\alpha}^{\text{el}} |E|}{2 D_{\alpha}} = \frac{h \mu_{\alpha}^{\text{el}} |E|}{2 \frac{\mu_{\alpha}^{\text{el}} V_{th}}{|z_{\alpha}|}} = |z_{\alpha}| \frac{h|E|}{2 V_{th}} > 1$$

$$\implies |\Delta\psi| > \frac{2 V_{th}}{|z_{\alpha}|} \leftarrow \text{advection-dominated problem}$$



Since in cellular electrophysiology  $V_m \sim 90\text{mV}$  the voltage drop between discretisation steps is in general much lower than  $2V_m \sim 50\text{mV}$ , there are typically no issues (such as spurious oscillations) when computing for the numerical solution of the problem.

Example:  $\vec{E} \equiv \text{const.}$  (it is often a reasonable approx.)



$$E = - \frac{\partial \Psi}{\partial x} = - \frac{\psi^{\text{out}} - \psi^{\text{in}}}{t_m} = \frac{\psi^{\text{in}} - \psi^{\text{out}}}{t_m} = \frac{V_m}{t_m}$$

$$V = \mu_{\alpha}^{\text{el}} \frac{|z_{\alpha}|}{z_{\alpha}} E = \mu_{\alpha}^{\text{el}} \frac{|z_{\alpha}|}{z_{\alpha}} \frac{V_m}{t_m}$$

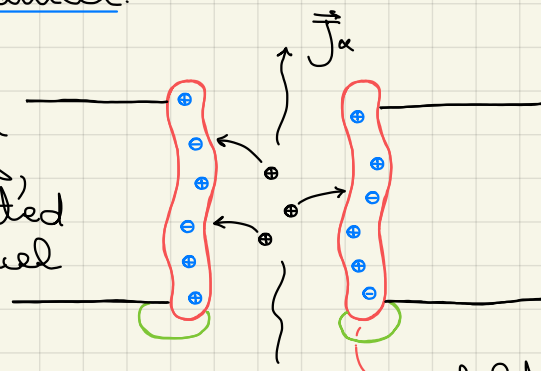
$$Pe_{\text{ad}} = \frac{h|V|}{2D_{\alpha}} = \frac{h \mu_{\alpha}^{\text{el}} |V_m|}{2 t_m \frac{\mu_{\alpha}^{\text{el}} V_m}{|z_{\alpha}|}} = |z_{\alpha}| \frac{h}{t_m} \frac{|V_m|}{2V_m} \approx 1.5$$

$Pe_{\text{ad}} < 1$

So as previously denoted we typically don't need any stabilization technique for the Péclet number.

This however can change in presence of fixed charges within the channel.

The electric field associated to these charges, which are located along the channel walls and should therefore be considered as boundary conditions, acts as a 2nd dimension that cannot be accounted for in our previous model.



can have fixed charges in their structure which generate an electric field

folded proteins specialized in the passage of specific ion species thanks to receptors

It is then clear that our 1D representation can be a good starting point for the study of this problem but it can definitely be improved.

Another aspect of higher detail that might play a significant role in the evaluation of the Péclet number is the intrachannel water.

As a matter of fact, so far we neglected any interaction between water molecules and ions (solvent and solute). To account for any possible "push-pull" mechanism between ions and fluid we use the following:

⇒ Velocity-extended Poisson-Nernst-Planck model:

$$\vec{J}_\alpha = q |z_\alpha| \mu_\alpha^{\text{el}} n_\alpha \vec{E} - q z_\alpha D_\alpha \vec{\nabla} n_\alpha + q z_\alpha n_\alpha \vec{U}_f$$

fluid velocity

It is a model for ion electro-fluid-diffusion.

$$\begin{aligned} \vec{f}_\alpha &= \frac{\vec{J}_\alpha}{q z_\alpha} = \frac{|z_\alpha|}{z_\alpha} \mu_\alpha^{\text{el}} n_\alpha \vec{E} + n_\alpha \vec{U}_f - D_\alpha \vec{\nabla} n_\alpha \\ &= n_\alpha \left( \vec{U}_f + \frac{|z_\alpha|}{z_\alpha} \mu_\alpha^{\text{el}} \vec{E} \right) - D_\alpha \vec{\nabla} n_\alpha \end{aligned}$$

✓ can now be greater than before!

Numerical oscillations due to  $Pe_{\text{ad}} > 1$  might now be relevant.

Consider now the simpler problem:

- Nernst-Planck model
- 1D case ( $x \in [0, t_m]$ )
- constant electric field ( $E = \frac{V_m}{t_m}$ )
- steady-state condition ( $\frac{\partial n_\alpha}{\partial t} = 0$ )

$$J_\alpha^{\text{cond}} = q \mu_\alpha^{\text{el}} |z_\alpha| n_\alpha \frac{V_m}{t_m} - q z_\alpha D_\alpha \frac{\partial n_\alpha}{\partial x}$$

The question is: how much is  $J_\alpha^{\text{cond}}$ ? and so  $gm$ ?

To solve this problem, we add a fifth hypothesis:

- $J_\alpha^{\text{cond}} = \text{const.}$  i.e. no consumption or generation inside the channel

$$\Rightarrow \frac{\partial J_\alpha^{\text{cond}}}{\partial x} = 0 \rightarrow \frac{\partial}{\partial x} \left( q \mu_\alpha^{\text{el}} |z_\alpha| n_\alpha \frac{V_m}{t_m} - q z_\alpha D_\alpha \frac{\partial n_\alpha}{\partial x} \right) = 0$$

$$\frac{\partial n_\alpha}{\partial x} \frac{V_m}{t_m} - \frac{V_{th}}{z_\alpha} \frac{\partial^2 n_\alpha}{\partial x^2} = 0$$

$$\begin{cases} -\frac{\partial^2 n_\alpha}{\partial x^2} + \frac{z_\alpha V_m}{V_{th} t_m} \frac{\partial n_\alpha}{\partial x} = 0, & x \in (0, t_m) \\ n_\alpha(0) = n_\alpha^{\text{in}} & n_\alpha(t_m) = n_\alpha^{\text{out}} \end{cases}$$

$n_\alpha(x)$  is the only unknown of the system

Characteristic polynomial  $P(\lambda)$

Eigenvalue analysis:  $-\lambda^2 + \frac{z_\alpha V_m}{V_{th} t_m} \lambda = 0$

$$\lambda \left( \frac{z_\alpha V_m}{V_{th} t_m} - \lambda \right) = 0$$

$$\lambda_1 = 0 \quad \lambda_2 = \gamma = \frac{z_\alpha V_m}{V_{th} t_m} \quad \lambda = [m^{-1}]$$

$$\Rightarrow n_\alpha(x) = A + B e^{\gamma x}$$

$$\begin{cases} A + B = n_\alpha^{\text{in}} \\ A + B e^{\gamma t_m} = n_\alpha^{\text{out}} \end{cases}$$

$$\begin{cases} A = \frac{n_\alpha^{\text{in}} e^{\gamma t_m} - n_\alpha^{\text{out}}}{e^{\gamma t_m} - 1} \\ B = \frac{n_\alpha^{\text{out}} - n_\alpha^{\text{in}}}{e^{\gamma t_m} - 1} \end{cases}$$

$$\Rightarrow J_\alpha^{\text{cond}} = q \mu_\alpha^{\text{el}} |z_\alpha| \frac{V_m}{t_m} (A + B e^{\gamma x}) - q z_\alpha D_\alpha B \gamma e^{\gamma x} = \text{const.}$$

$$= q \mu_\alpha^{\text{el}} |z_\alpha| \frac{V_m}{t_m} A + q \mu_\alpha^{\text{el}} |z_\alpha| \frac{V_m}{t_m} B e^{\gamma x} - q z_\alpha \frac{\mu_\alpha^{\text{el}} V_{th}}{|z_\alpha|} B \frac{z_\alpha V_m}{V_{th} t_m} e^{\gamma x}$$

correctly does not depend on  $x$  as it was supposed to be constant

equal and opposite

$$= q \mu_\alpha^{\text{el}} |z_\alpha| \frac{V_m}{t_m} \frac{n_\alpha^{\text{in}} e^{\gamma t_m} - n_\alpha^{\text{out}}}{e^{\gamma t_m} - 1}$$

$$\beta_m = z_\alpha \frac{V_m}{V_{th}} = \gamma t_m$$

$$= q \frac{D_\alpha}{t_m} z_\alpha \beta_m \frac{n_\alpha^{\text{in}} e^{\beta_m} - n_\alpha^{\text{out}}}{e^{\beta_m} - 1}$$

Introducing the inverse of the Bernoulli function:

$$\left[ \text{Be}(x) := \frac{x}{e^x - 1} \right]$$

we can then write:  $J_{\alpha}^{\text{cond}} = q \frac{D_{\alpha}}{t_m} z_{\alpha} \text{Be}(\beta_m) [n_{\alpha}^{\text{in}} e^{\beta_m} - n_{\alpha}^{\text{out}}]$

$$\text{Be}(x) \cdot e^x = \frac{x}{1 - e^{-x}} = \frac{-x}{e^{-x} - 1} = \text{Be}(-x)$$

$$= q \frac{D_{\alpha}}{t_m} z_{\alpha} [\text{Be}(-\beta_m) n_{\alpha}^{\text{in}} - \text{Be}(\beta_m) n_{\alpha}^{\text{out}}]$$

2) Goldman - Hodgkin - Katz (GHK) model:

$$J_{\alpha}^{\text{GHK}} = -q \frac{D_{\alpha}}{t_m} z_{\alpha} [\text{Be}(\beta_m) n_{\alpha}^{\text{out}} - \text{Be}(-\beta_m) n_{\alpha}^{\text{in}}] = \frac{I_{\alpha}}{S_{\text{cell}}}$$

with  $[\beta_m = z_{\alpha} \frac{V_m}{V_{th}}]$

$$g_{m,\alpha} = q \frac{z_{\alpha}^2 n_{\alpha}^{\text{in}}}{V_{th}} \frac{D_{\alpha}}{t_m} \quad \text{effective conductance} \quad [\frac{\Omega^{-1}}{m^2}]$$

$$E_{m,\alpha} = \frac{V_{th}}{z_{\alpha}} \left( \frac{n_{\alpha}^{\text{out}}}{n_{\alpha}^{\text{in}}} - 1 \right) \text{Be}(\beta_m) \quad \text{effective Nernst potential} \quad [V]$$

In general it is very hard to know the exact value of the diffusivity  $D_{\alpha}$  within the membrane channel, as much as it is hard to know its thickness  $t_m$ . However, it is possible to experimentally derive the membrane permeability (with respect to a certain ion species  $\alpha$ ):

$$P_{\alpha} = \frac{D_{\alpha}}{t_m} = \left[ \frac{m}{s} \right]$$

Note: the GHK result for conduction current resembles the general formula for current density:

$$\vec{J}_{\alpha} = q n_{\alpha} z_{\alpha} \vec{U}_{\alpha} \quad \rightarrow \text{drift velocity}$$

It is then evident that the GHK current displays the same terms through the use of effective parameters:

$$J_{\alpha}^{\text{GHK}} = q n_{\alpha}^{\text{eff}} z_{\alpha} P_{\alpha}$$

where  $n_{\alpha}^{\text{eff}} = \text{Be}(-\beta_m) n_{\alpha}^{\text{in}} - \text{Be}(\beta_m) n_{\alpha}^{\text{out}}$  and  $P_{\alpha} = \frac{D_{\alpha}}{t_m}$  are the effective parameters for concentration and velocity, in the sense that they efficiently approximate with a single constant term their respective variable, which in reality may vary



in both time and space along the membrane channel.

The whole GHK equation is an efficient model to approximate with a constant effective current the actual channel current.

Let's now see how well this model can estimate the conduction current in some limit cases.

$$J_{\alpha}^{\text{GHK}} = -q P_{\alpha} z_{\alpha} \left[ \text{Be}(\beta_m) n_{\alpha}^{\text{out}} - \text{Be}(-\beta_m) n_{\alpha}^{\text{in}} \right]$$

VS.

$$J_{\alpha}^{\text{cond}} = q \mu_{\alpha}^{\text{el}} |z_{\alpha}| n_{\alpha} \frac{V_m}{t_m} - q z_{\alpha} D_{\alpha} \frac{\partial n_{\alpha}}{\partial x}$$

$$\text{Be}(0) = 1$$

1. Zero electric field  $\rightarrow E = 0, V_m = 0, \beta_m = 0$

$$J_{\alpha}^{\text{GHK}} = -q z_{\alpha} \frac{D_{\alpha}}{t_m} (n_{\alpha}^{\text{out}} - n_{\alpha}^{\text{in}}) \quad \text{VS.} \quad J_{\alpha}^{\text{cond}} = -q z_{\alpha} D_{\alpha} \frac{\partial n_{\alpha}}{\partial x}$$

The GHK formula approximates the concentration gradient through the incremental ratio between the two endpoints of the domain:

$$\frac{\partial n_{\alpha}}{\partial x} \approx \frac{n_{\alpha}^{\text{out}} - n_{\alpha}^{\text{in}}}{t_m}$$

2. Uniform concentration  $\rightarrow n_{\alpha}^{\text{out}} = n_{\alpha}^{\text{in}} = \bar{n}_{\alpha}$

$$J_{\alpha}^{\text{GHK}} = -q z_{\alpha} \frac{D_{\alpha}}{t_m} \bar{n}_{\alpha} \left[ \text{Be}(\beta_m) - \text{Be}(-\beta_m) \right] =$$

$$x + \text{Be}(x) = x + \frac{x}{e^x - 1} = \frac{x e^x}{e^x - 1} = e^x \text{Be}(x) = \text{Be}(-x)$$

$$= -q z_{\alpha} \frac{D_{\alpha}}{t_m} \bar{n}_{\alpha} [-\beta_m]$$

$$= q z_{\alpha} \frac{1}{t_m} \left( \mu_{\alpha}^{\text{el}} \frac{V_{\text{th}}}{|z_{\alpha}|} \right) \bar{n}_{\alpha} \left( z_{\alpha} \frac{V_m}{V_{\text{th}}} \right)$$

$$= q |z_{\alpha}| \mu_{\alpha}^{\text{el}} \bar{n}_{\alpha} \frac{V_m}{t_m} \quad \text{VS.} \quad J_{\alpha}^{\text{cond}} = q |z_{\alpha}| \mu_{\alpha}^{\text{el}} \bar{n}_{\alpha} \frac{V_m}{t_m}$$

The GHK formula is exactly equal to the theoretical one, which is reasonable since  $J_{\alpha}^{\text{cond}}$  is constant and  $J_{\alpha}^{\text{GHK}}$  is also constant by definition, so it has no issues with approximating the real behaviour.

As the GHK approximation is reliable in both limit cases, it is expected to always yield a reliable value for any intermediate case.

We will hereafter consider the following velocity-extended problem:

- Poisson-Nernst-Planck model
- 3D case

to evaluate the complete expression of the drift velocity  $\vec{U}_\alpha$ .

$$\vec{J}_\alpha = q |z_\alpha| \mu_\alpha^{\text{el}} n_\alpha \vec{E} - q z_\alpha D_\alpha \vec{\nabla} n_\alpha + q z_\alpha n_\alpha \vec{U}_f \iff \vec{J}_\alpha = q n_\alpha z_\alpha \vec{U}_\alpha$$

$$\vec{U}_\alpha = ?$$

$$\begin{aligned} \vec{J}_\alpha &= q |z_\alpha| \mu_\alpha^{\text{el}} n_\alpha \vec{E} - q z_\alpha D_\alpha \vec{\nabla} n_\alpha + q z_\alpha n_\alpha \vec{U}_f \\ &= q z_\alpha n_\alpha \left[ \vec{U}_f - \mu_\alpha^{\text{el}} \frac{|z_\alpha|}{z_\alpha} \vec{\nabla} \psi - D_\alpha \frac{1}{n_\alpha} \vec{\nabla} n_\alpha \right] \\ &= q z_\alpha n_\alpha \left[ \vec{U}_f - \mu_\alpha^{\text{el}} \frac{|z_\alpha|}{z_\alpha} \vec{\nabla} \psi - \mu_\alpha^{\text{el}} \frac{V_{\text{th}}}{|z_\alpha|} \vec{\nabla} \ln \left( \frac{n_\alpha}{n_{\text{ref}}} \right) \right] \\ &= q z_\alpha n_\alpha \left[ \vec{U}_f - \frac{\mu_\alpha^{\text{el}} V_{\text{th}}}{|z_\alpha|} \left( \vec{\nabla} \ln \left( \frac{n_\alpha}{n_{\text{ref}}} \right) + \frac{|z_\alpha|^2}{z_\alpha} \frac{1}{V_{\text{th}}} \vec{\nabla} \psi \right) \right] \\ &= q z_\alpha n_\alpha \left[ \vec{U}_f - \frac{\mu_\alpha^{\text{el}} V_{\text{th}}}{|z_\alpha|} \frac{z_\alpha}{V_{\text{th}}} \vec{\nabla} \left( \psi + \frac{V_{\text{th}}}{z_\alpha} \ln \left( \frac{n_\alpha}{n_{\text{ref}}} \right) \right) \right] \end{aligned}$$

$\downarrow$  electric potentials
 $\downarrow$  chemical potentials

Electrochemical potential  $\left[ \varphi_\alpha^{\text{EC}} := \psi + \frac{V_{\text{th}}}{z_\alpha} \ln \left( \frac{n_\alpha}{n_{\text{ref}}} \right) \right]$

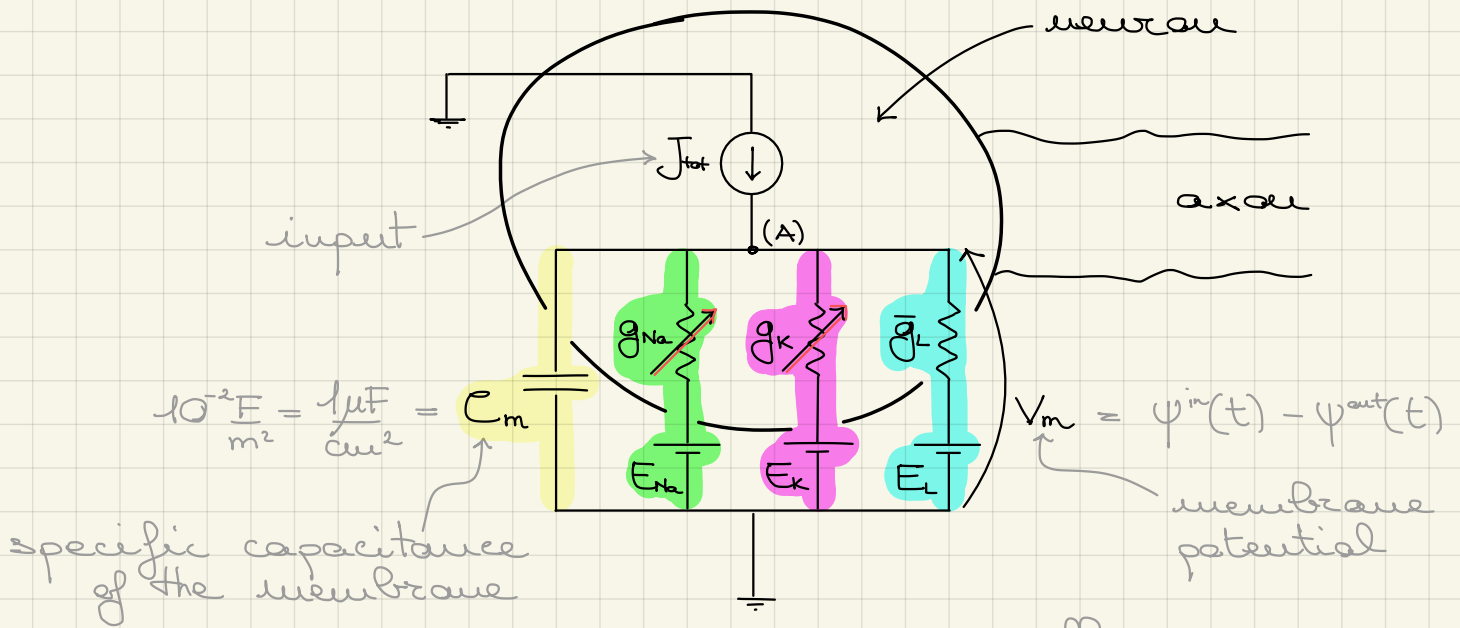
$$\vec{J}_\alpha = q z_\alpha n_\alpha \left[ \vec{U}_f - \mu_\alpha^{\text{el}} \frac{z_\alpha}{|z_\alpha|} \vec{\nabla} \varphi_\alpha^{\text{EC}} \right]$$

electrochemical field  $\left[ \vec{E}_\alpha^{\text{EC}} = -\vec{\nabla} \varphi_\alpha^{\text{EC}} \right]$

$$\vec{J}_\alpha = q n_\alpha z_\alpha \vec{U}_\alpha \longrightarrow \vec{U}_\alpha = \vec{U}_f + \mu_\alpha^{\text{el}} \frac{z_\alpha}{|z_\alpha|} \vec{E}_\alpha^{\text{EC}} \longrightarrow \begin{cases} \vec{E}_\alpha^{\text{EC}} = -\vec{\nabla} \varphi_\alpha^{\text{EC}} \\ \varphi_\alpha^{\text{EC}} = \psi + \frac{V_{\text{th}}}{z_\alpha} \ln \left( \frac{n_\alpha}{n_{\text{ref}}} \right) \end{cases}$$



### 3) Hodgkin - Huxley (neuron) cell model and characterization of the gating variables

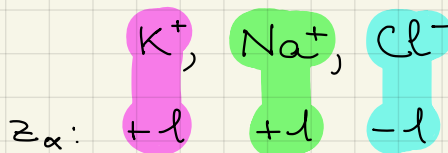


$\bar{g}_L$ : "leakage conductance" } associated with chloride ions  
 $E_L$ : Nernst potential } (+ some other ions)

$g_K = g_K(t, V_m)$ : conductance } associated with potassium ions  
 $E_K$ : Nernst potential }

$g_{Na} = g_{Na}(t, V_m)$ : conductance } associated with sodium ions  
 $E_{Na}$ : Nernst potential }

variable parameters



KCL at node (A):

$$-J_{tot} + C_m \frac{dV_m}{dt} + g_{Na}(V_m - E_{Na}) + g_K(V_m - E_K) + \bar{g}_L(V_m - E_L) = 0$$

with initial conditions:

$$V_m(t = t_0) = V_m^i$$

If we manage to find an expression for  $g_{Na}$  and  $g_K$  then we know how  $V_m$  will evolve over time.

Notation change:  $J_\alpha = g_\alpha \cdot (V_m - E_{m,\alpha})$   
 $= g_\alpha \cdot (V_m - E_m + E_m - E_{m,\alpha})$

where  $E_m$  is the resting (Nernst) potential of the cell while  $E_{m,\alpha}$  is the Nernst potential associated to ion species  $\alpha$  (e.g.  $E_{Na}$ ).  
 $E_m$  is a constant that holds information related to each  $E_{m,\alpha}$ .

$$\rightarrow \boxed{J_\alpha = g_\alpha \cdot (V - V_\alpha)}$$

where  $V := V_m - E_m$  and  $V_\alpha = E_{m,\alpha} - E_m$

With this change in the notation it is easier to see when the cell is not at equilibrium ( $V \neq 0$ ).

$$\Rightarrow \boxed{-J_{tot} + C_m \frac{dV}{dt} + g_{Na}(V - V_{Na}) + g_K(V - V_K) + \bar{g}_L(V - V_L) = 0}$$

$E_m = \text{const.}$

$t_0 = 0$

equilibrium at starting condition

Nernst potential for ion species  $\alpha$ :  $E_{m,\alpha} = \frac{V_{th}}{z_\alpha} \ln\left(\frac{C_\alpha^{out}}{C_\alpha^{in}}\right)$

Nernst potential of the cell:  $E_m = ?$

$$\updownarrow$$

$$J_\alpha^{cond} = 0$$

Assumption: all ions are monovalent ( $z_\alpha = \pm 1$ )

Goldman potential:

$$E_m := V_{th} \ln \left[ \frac{\sum_{i=1}^{M^+} P_i^+ C_i^{+(out)} + \sum_{i=1}^{M^-} P_i^- C_i^{-(in)}}{\sum_{i=1}^{M^+} P_i^+ C_i^{+(in)} + \sum_{i=1}^{M^-} P_i^- C_i^{-(out)}} \right]$$

$\updownarrow$   
 $J_{tot}^{cond} = 0$

permeability  $\rightarrow$   
concentration  $\rightarrow$

where  $J_{tot}^{cond} = \sum_{i=1}^{M^+} J_i^+ + \sum_{i=1}^{M^-} J_i^-$

$M^+ = 2$  number of cationic species ( $K^+$ ,  $Na^+$ )

$M^- = 1$  number of anionic species ( $Cl^-$ )

## Gating variables

$g_{\alpha} = g_{\alpha}(s_{\alpha})$  where  $s_{\alpha}$  are "gating variables"

$$s_{\alpha} = s_{\alpha}(t, v) \quad 0 \leq s_{\alpha i} \leq 1$$

$\swarrow$   $\alpha$ -ion channel closed       $\searrow$   $\alpha$ -ion channel open

We now need an expression for  $g_{\alpha}(s_{\alpha})$  and, more importantly, for  $s_{\alpha}(t, v)$ .

From Hodgkin-Huxley (1952): (i) Potassium

$$g_K = \bar{g}_K \cdot n^4 \rightarrow S_g = n$$

differential equation  
(Balance equation)

$$\frac{dn}{dt} = \underbrace{\alpha_n(1-n)}_{\text{generation rate}} - \underbrace{\beta_n n}_{\text{consumption rate}}$$

$$\alpha_n = (0,1 \cdot \frac{1}{\text{ms}}) \cdot \text{Be}\left(\frac{-v}{10\text{mV}} + 1\right)$$

$$\beta_n = (0,125 \cdot \frac{1}{\text{ms}}) e^{-v/80\text{mV}}$$

generation rate      consumption rate

(ii) Sodium

$$g_{Na} = \bar{g}_{Na} m^3 h \rightarrow \vec{S}_g = [m, h]$$

$$\frac{dm}{dt} = \alpha_m(1-m) - \beta_m m$$

$$\frac{dh}{dt} = \alpha_h(1-h) - \beta_h h$$

$$\alpha_m = \left(\frac{1}{\text{ms}}\right) \cdot \text{Be}\left(\frac{-v}{10\text{mV}} + 2,5\right)$$

$$\beta_m = \left(4 \cdot \frac{1}{\text{ms}}\right) e^{-v/48\text{mV}}$$

$$\alpha_h = \left(0,07 \cdot \frac{1}{\text{ms}}\right) e^{-v/20\text{mV}}$$

$$\beta_h = \frac{\left(\frac{1}{\text{ms}}\right)}{e^{(v/10\text{mV} + 3)} + 1}$$

$$\bar{g}_{\alpha} = \text{const.} = \left[\frac{\Omega^{-1}}{\text{m}^2}\right] \quad \alpha_{S_g} = \alpha_{S_g}(v) = [\text{s}^{-1}] \quad \beta_{S_g} = \beta_{S_g}(v) = [\text{s}^{-1}]$$

$S_g$ : proportion of ions inside the membrane

$1 - S_g$ : " " " outside " "

$\alpha_{S_g}$ : rate of transfer from outside to inside

$\beta_{S_g}$ : " " " inside " outside

These equations were the result of a careful interpolation between mathematical models and experimental data.

We finally have everything we needed to determine the membrane conductances and therefore the solution of the Hodgkin-Huxley model:

$$C_m \frac{dV}{dt} = J_{tot} - J_{Na} - J_K - J_L$$

$$J_{Na} = g_{Na}(V - V_{Na})$$

$$J_K = g_K(V - V_K)$$

$$J_L = \bar{g}_L(V - V_L)$$

$$g_K = \bar{g}_K n^4$$

$$g_{Na} = \bar{g}_{Na} m^3 h$$

$$\frac{dn}{dt} = \alpha_n(V)(1-n) - \beta_n(V)n$$

$$\frac{dm}{dt} = \alpha_m(V)(1-m) - \beta_m(V)m$$

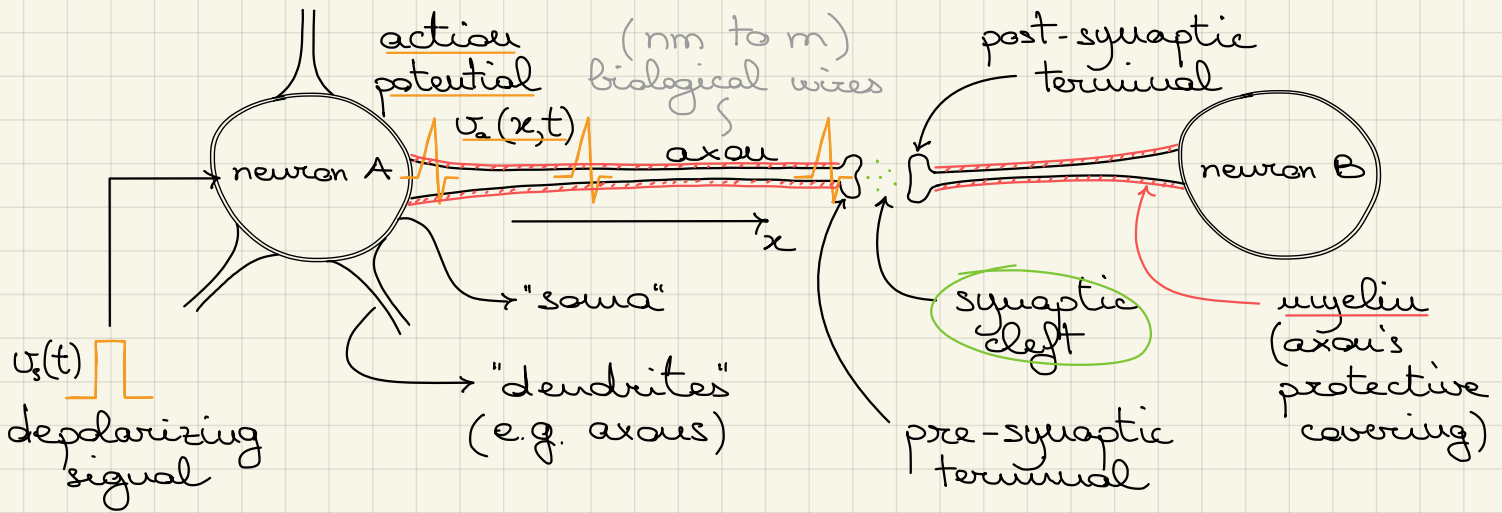
$$\frac{dh}{dt} = \alpha_h(V)(1-h) - \beta_h(V)h$$

$$\begin{cases} \frac{dy}{dt} = f(t, y(t)) \\ y(t_0) = y_0 \end{cases}$$

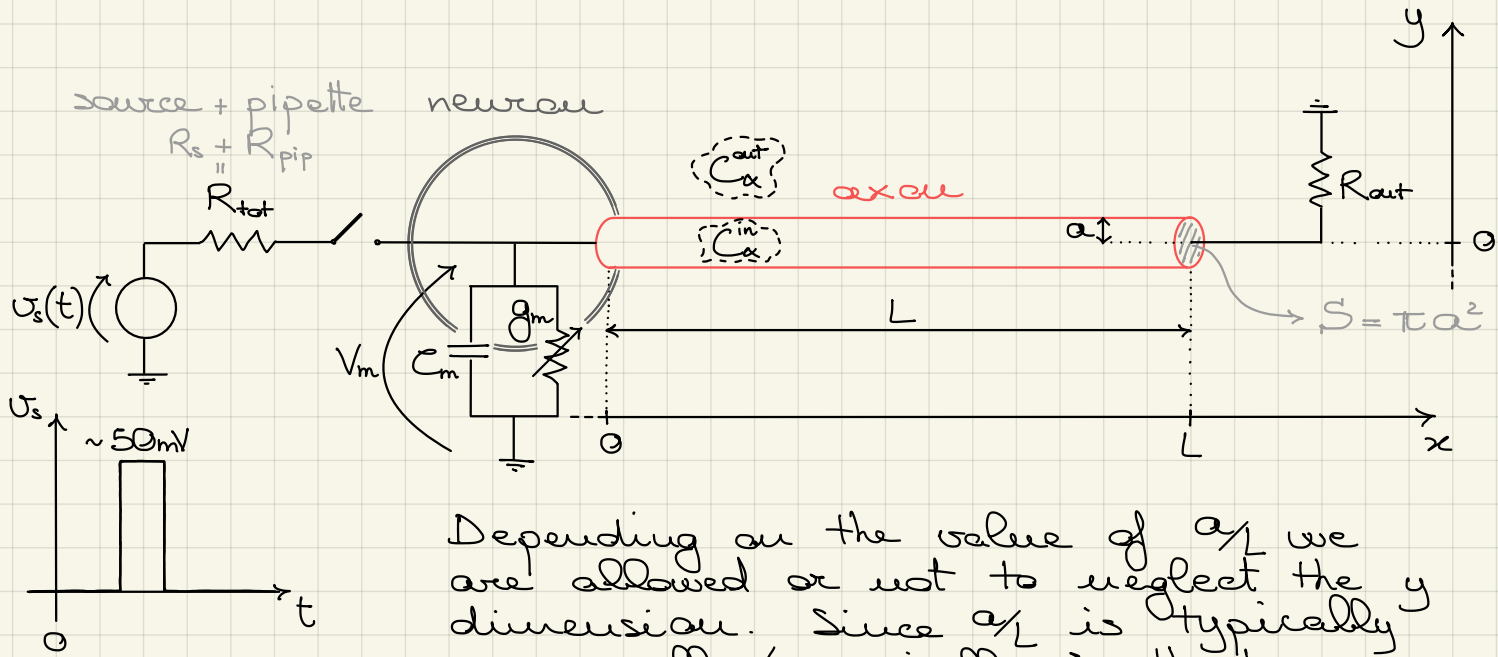
$$y = [V, n, m, h]^T$$

non-linear  
Cauchy problem!  
(use ODE15s)

# The Cable Equation model



It's not just a time-varying problem but it is space-varying too!

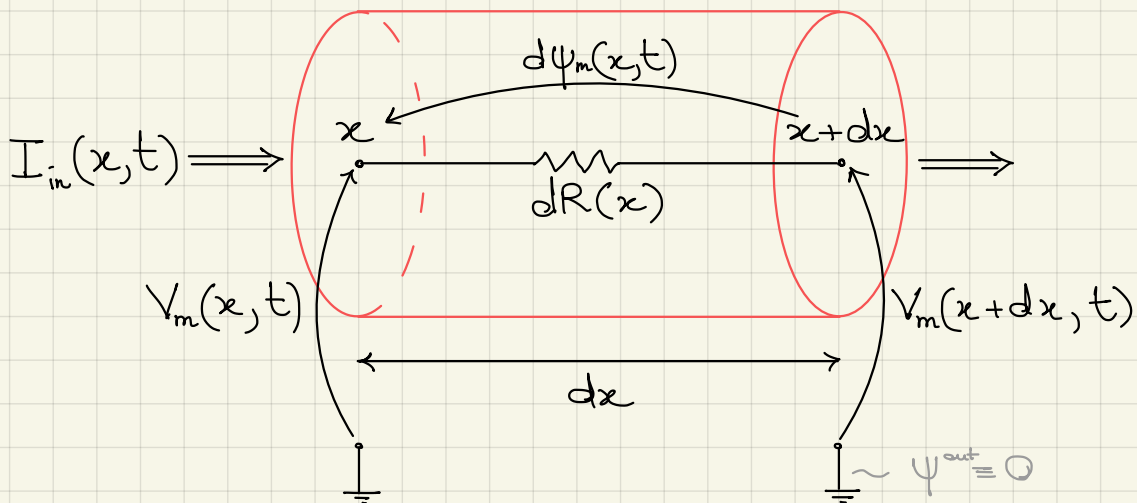
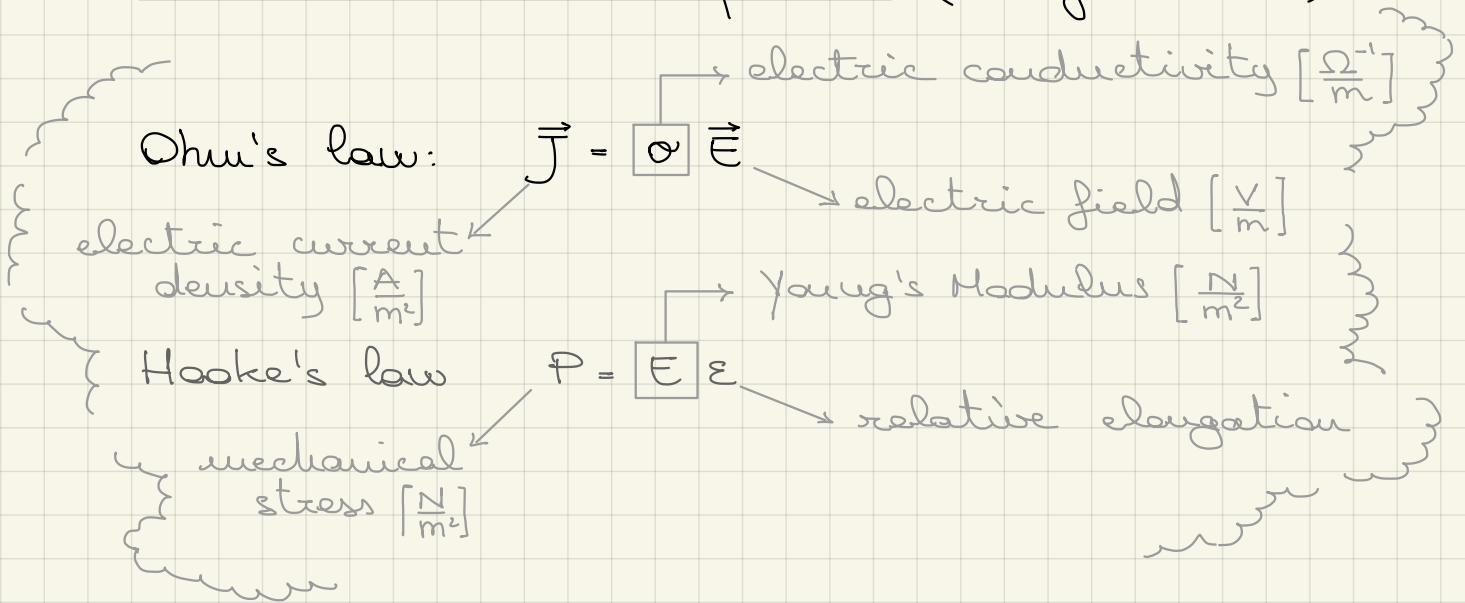


Depending on the value of  $a/L$  we are allowed or not to neglect the  $y$  dimension. Since  $a/L$  is typically very small (especially in the human body:  $a/L \approx 5 \cdot 10^{-5}$ ) we can consider just the  $x$  dimension in our problem.

- Assumptions:
1. in the intra/extracellular region  $C_x = C_x(x,t) = \text{const.}$
  2.  $\psi^{\text{out}}(x,t) = 0$  → concentration
  3. in the intracellular region  $\rho_{ax} = \text{const.} > 0$
- $\left[ \frac{\Omega}{m} \right]$  resistivity

# Characterization of the axon duct

## 1) Current constitutive equation (along x axis)



$$dR(x) = \frac{\rho_{ax} dx}{S} = \frac{\rho_{ax} dx}{\pi a^2}$$

$$dV_m(x,t) = V_m(x,t) - V_m(x+dx,t)$$

$$= dR(x) \cdot I_{in}(x,t)$$

$$= \frac{\rho_{ax} dx}{\pi a^2} \cdot I_{in}(x,t)$$

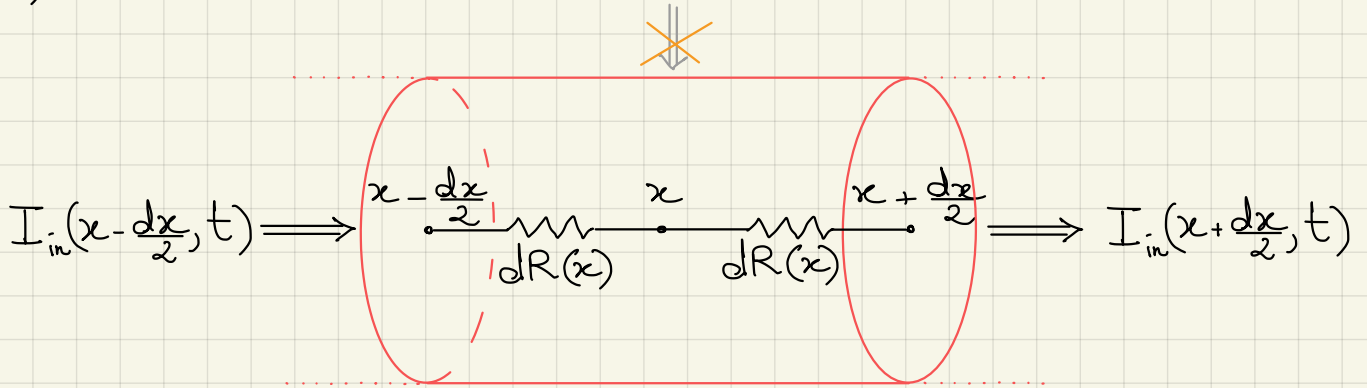
$$I_{in}(x,t) = \frac{dV_m(x,t)}{dx} \frac{\pi a^2}{\rho_{ax}}$$

$$= - \frac{V_m(x+dx,t) - V_m(x,t)}{dx} \frac{\pi a^2}{\rho_{ax}}$$

$$dx \rightarrow 0 \longrightarrow \boxed{I_{in}(x,t) = - \frac{\pi a^2}{\rho_{ax}} \frac{\partial V_m(x,t)}{\partial x}}$$

$$\text{Ohm's law: } J_x = \sigma E_x = \frac{1}{S} \left( - \frac{\partial \Psi_x}{\partial x} \right)$$

## 2) Current balance law

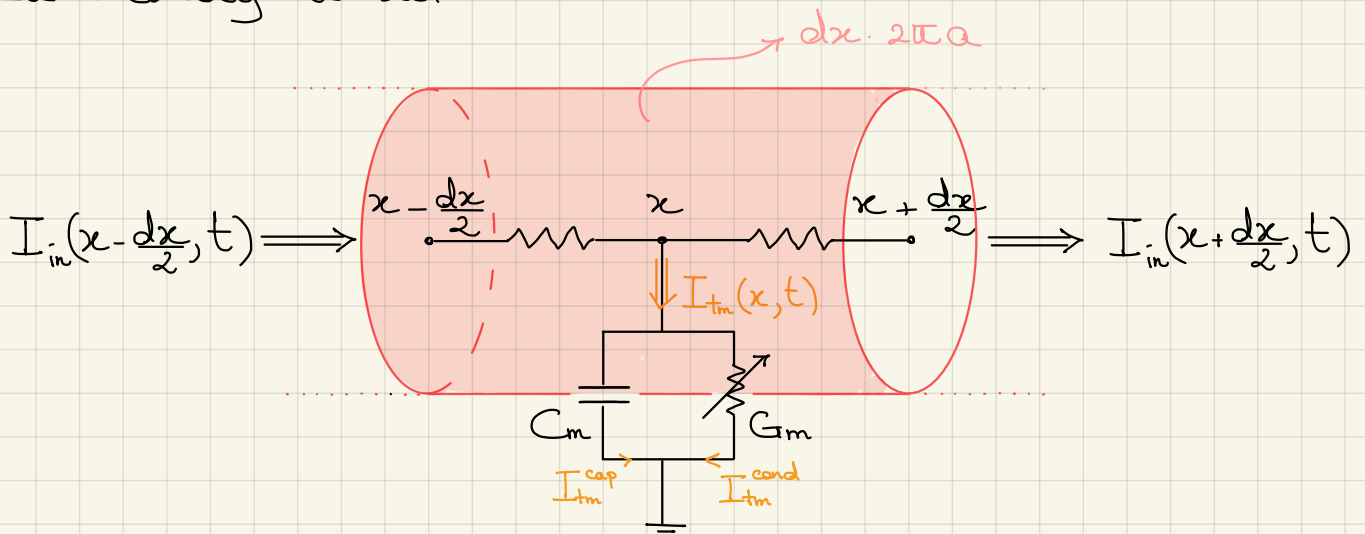


$$\frac{I_{in}(x + \frac{dx}{2}, t) - I_{in}(x - \frac{dx}{2}, t)}{dx} = 0$$

$$dx \rightarrow 0 \longrightarrow \boxed{\frac{\partial I_{in}(x, t)}{\partial x} = 0} \quad \leftarrow \text{no current in the vertical direction}$$

$$\begin{cases} 1) \\ 2) \end{cases} \Longrightarrow \boxed{-\frac{\partial^2 V_m(x, t)}{\partial x^2} = 0} \quad \leftarrow \text{Laplace equation}$$

However this assumption was not correct. In reality it is:



$$\text{KCL at } x: I_{in}(x + \frac{dx}{2}, t) - I_{in}(x - \frac{dx}{2}, t) + I_{tm}(x, t) = 0$$

$$I_{in}(x + \frac{dx}{2}, t) - I_{in}(x - \frac{dx}{2}, t) + [I_{tm}^{cond}(x, t) + I_{tm}^{cap}(x, t)] = 0$$

$$* I_{tm}^{cap}(x, t) = j_{tm}^{cap}(t, m) \cdot dx \cdot 2\pi a = \epsilon_m \frac{\partial V_m(x, t)}{\partial t} dx \cdot 2\pi a$$



Assumption: use the linear resistor model to characterize the conductive  $t_m$  current

$$\# I_{t_m}^{\text{cond}}(x,t) = \sum_{\alpha} j_{t_m,\alpha}^{\text{cond}}(x,t) \cdot dx \cdot 2\pi a =$$

$$= \left[ \sum_{\alpha} g_{m,\alpha} (V_m(x,t) - E_{m,\alpha}) \right] dx \cdot 2\pi a$$

$$\longrightarrow \left[ \sum_{\alpha} g_{m,\alpha} (V_m(x,t) - E_{m,\alpha}) + c_m \frac{\partial V_m(x,t)}{\partial t} \right] dx \cdot 2\pi a =$$

$$= \frac{I_{in}(x - \frac{dx}{2}, t) - I_{in}(x + \frac{dx}{2}, t)}{dx \cdot 2\pi a}$$

$$\longrightarrow \left\{ \begin{array}{l} c_m \frac{\partial V_m(x,t)}{\partial t} + \frac{1}{2\pi a} \frac{\partial I_{in}(x,t)}{\partial x} + \left( \sum_{\alpha} g_{m,\alpha} \right) V_m(x,t) = \sum_{\alpha} g_{m,\alpha} E_{m,\alpha} \\ I_{in}(x,t) = -\frac{\pi a^2}{\int_{ax}} \frac{\partial V_m(x,t)}{\partial x} \\ \text{i.e. } V_m(x,0) = V_m^0(x) \quad \forall x \in (0,L) \\ \text{B.c. for } V_m \text{ at } x=0, x=L \quad \forall t \in (0, T) \end{array} \right.$$

Linear Cable eq. model

it's a PARABOLIC problem (like heat equation)

Remove the previous linear resistor model assump.

Non-linear resistor model:  $j_{t_m,\alpha}^{\text{cond}} = g_{m,\alpha}(V_m) [V_m - E_{m,\alpha}(V_m)]$

$$\longrightarrow \left\{ \begin{array}{l} c_m \frac{\partial V_m(x,t)}{\partial t} + \frac{1}{2\pi a} \frac{\partial I_{in}(x,t)}{\partial x} + \left[ \sum_{\alpha} g_{m,\alpha}(V_m) \right] V_m(x,t) = \sum_{\alpha} g_{m,\alpha}(V_m) E_{m,\alpha}(V_m) \\ I_{in}(x,t) = -\frac{\pi a^2}{\int_{ax}} \frac{\partial V_m(x,t)}{\partial x} \\ \text{i.e. } V_m(x,0) = V_m^0(x) \quad \forall x \in (0,L) \\ \text{B.c. for } V_m \text{ at } x=0, x=L \quad \forall t \in (0, T) \end{array} \right.$$

Non-linear Cable eq. model

Note the familiar form of the problem:

$$(1) \frac{\partial u}{\partial t} + \frac{\partial J}{\partial x} = g - k u \quad (2) J = -D \frac{\partial u}{\partial x} + V u$$

$G - R$

$$(1) \frac{\partial V_m}{\partial t} + \frac{1}{2\pi a \epsilon_m} \frac{\partial I_{in}(V_m)}{\partial x} = - \sum_{\alpha} \frac{g_{m,\alpha}(V_m)}{\epsilon_m} [V_m - E_{m,\alpha}(V_m)]$$

$$= \underbrace{\sum_{\alpha} \frac{g_{m,\alpha}(V_m) E_{m,\alpha}(V_m)}{\epsilon_m}}_{G(V_m)} - \underbrace{\frac{g_{tot}(V_m)}{\epsilon_m} \cdot V_m}_{R(V_m)}$$

$$(2) I_{in}(V_m) = - \frac{\pi a^2}{g_{ax}} \frac{\partial V_m}{\partial x}$$

$$u = V_m \quad J(u) = I_{in}(V_m) \quad D = \frac{a}{2\epsilon_m g_{ax}} \quad V = 0$$

$$g(u) = \sum_{\alpha} \frac{g_{m,\alpha}(V_m) E_{m,\alpha}(V_m)}{\epsilon_m} \quad K(u) = \frac{g_{tot}(V_m)}{\epsilon_m}$$

Since this problem involves both time and space variables, as already seen we must proceed to time semidiscretization

↳ θ-method with  $\theta = 1$

(i.e. Backward Euler method, only first order converg but unconditionally stable + positivity preserving)



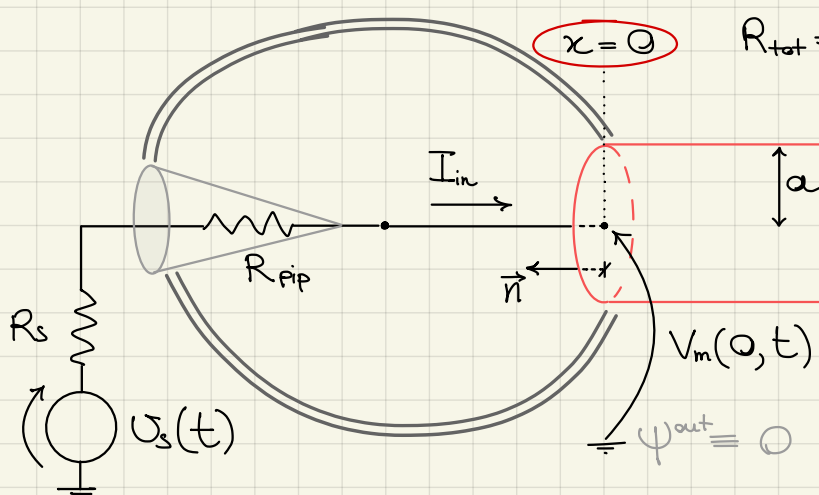
$$N_T \geq 1 \quad \Delta t = \frac{T}{N_T} \quad n = 0, 1, \dots, N_T$$

$$\rightarrow \frac{V_m^{n+1}}{\Delta t} + \frac{1}{2\pi a \epsilon_m} \frac{\partial I_{in}(V_m^{n+1})}{\partial x} = \frac{V_m^n}{\Delta t} + \frac{1}{\epsilon_m} \left[ \sum_{\alpha} g_{m,\alpha}(V_m^{n+1}) E_{m,\alpha}(V_m^{n+1}) - g_{tot}(V_m^{n+1}) \cdot V_m^{n+1} \right]$$

given i.c.:  $V_m^0 = V_m(x, 0)$

and b.c.: ?

### Boundary conditions



$$R_{tot} = R_s + R_{pip} \rightarrow U_s(t) - R_{tot} I_{in}(0, t) - V_m(0, t) = 0$$

$$- I_{in}(0, t) = \frac{1}{R_{tot}} V_m(0, t) - \frac{U_s(t)}{R_{tot}}$$

Robin b.c.

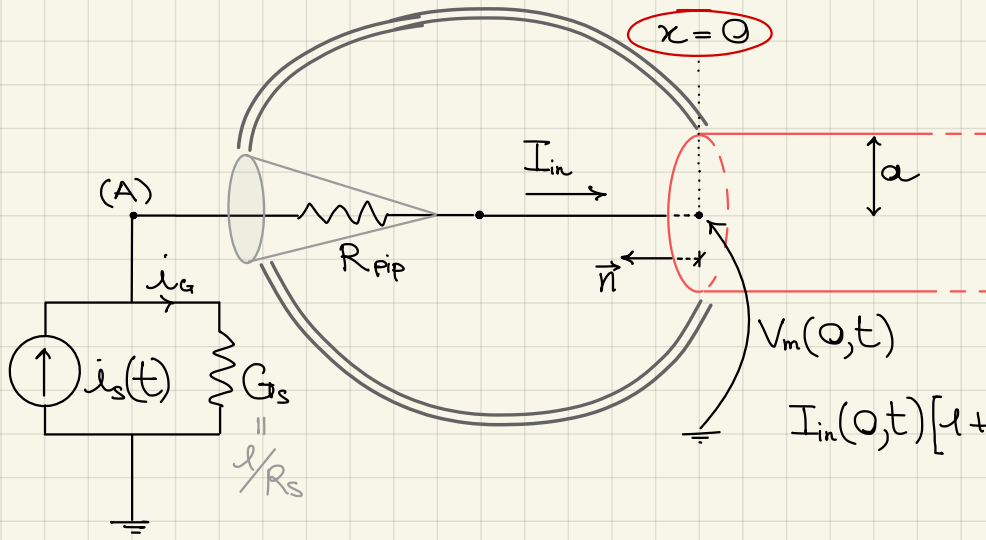
$$g \vec{j} \cdot \vec{n} = \alpha u - \beta$$

$$\gamma_0 \vec{I}_{in}(0,t) \cdot \vec{n}|_{x=0} = \alpha_0 V_m(0,t) - \beta_0(t)$$

voltage clamp boundary condition

where  $\gamma_0 = 1$     $\alpha_0 = \frac{1}{R_{pip}}$     $\beta_0 = \frac{V_s(t)}{R_{tot}}$

What about a current clamp experiment (i.e. current source)?



$$\psi_A(t) - R_{pip} I_{in}(0,t) - V_m(0,t) = 0$$

$$i_s(t) - I_{in}(0,t) + i_G(t) = 0$$

$$i_G(t) = G_s \psi_A(t)$$

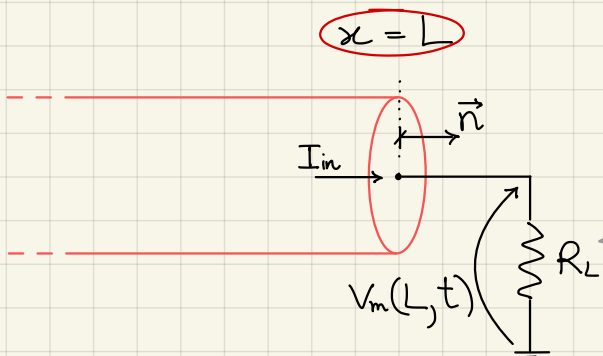
$$I_{in}(0,t) [1 + G_s R_{pip}] - i_s(t) + G_s V_m(0,t) = 0$$

$$-I_{in}(0,t) = \frac{G_s}{1 + G_s R_{pip}} V_m(0,t) - \frac{i_s(t)}{1 + G_s R_{pip}}$$

$$\gamma_0 \vec{I}_{in}(0,t) \cdot \vec{n}|_{x=0} = \alpha_0 V_m(0,t) - \beta_0(t)$$

current clamp boundary condition

where  $\gamma_0 = 1$     $\alpha_0 = \frac{G_s}{1 + G_s R_{pip}}$     $\beta_0 = \frac{i_s(t)}{1 + G_s R_{pip}}$



"load" resistance that approximates the synaptic cleft environment at the axon's termination

$$V_m(L,t) = R_L \cdot I_{in}(L,t) \implies I_{in}(L,t) = \frac{1}{R_L} V_m(L,t)$$

$$\gamma_L \vec{I}_{in}(L,t) \cdot \vec{n}|_{x=L} = \alpha_L V_m(L,t) - \beta_L(t)$$

synaptic end boundary condition

where  $\gamma_L = 1$     $\alpha_L = \frac{1}{R_L}$     $\beta_L = 0$

Our problem has, at this point, the following form:

given  $V_m = V_m(x)$ ,  $x \in \Omega = (0, L)$ ,  $\forall n = 0, 1, \dots, N_T$  solve

$$\begin{cases} \frac{1}{2\pi a \epsilon_m} \frac{\partial I_{in}(V_m^{n+1})}{\partial x} + \frac{V_m^{n+1}}{\Delta t} = \frac{V_m^n}{\Delta t} - \frac{1}{C_m} j_{tm}^{cond}(x, V_m^{n+1}, C_\alpha, S_g^{n+1}) \\ I_{in}(V_m) = -\frac{\pi a^2}{\rho_{ax}} \frac{\partial V_m}{\partial x} \rightarrow -\frac{\partial V}{\partial x} = E \text{ quasi-static approx.} \\ j_{tm}^{cond} = R - G = V_m \cdot K - G \\ \rho_{ax} \vec{I}_{in} \cdot \vec{n} \Big|_{\partial \Omega} = \alpha_{ax} V_m^{n+1} - \beta_{ax}^{n+1} \end{cases}$$

depends on the model  
(no current, linear resistor, GHK, H-H, etc.)

non-linear cable equation

"quasi-linear"

$$\rightarrow N(u) = 0$$

zero order

since the non-linearity is in the zero order terms

$$N(u) = \epsilon \frac{\partial I_{in}(u)}{\partial x} + \sigma u - g(u) + r(u)$$

$$u := V_m^{n+1} \quad \epsilon = \frac{1}{2\pi a \epsilon_m} \quad \sigma = \frac{1}{\Delta t}$$

$$G = g(u) = \frac{1}{C_m} \sum_{\alpha} g_{\alpha}(u) E_{\alpha}(u) \quad r(u) = -\frac{V_m^n}{\Delta t} + \frac{u}{C_m} g_{tot}(u)$$

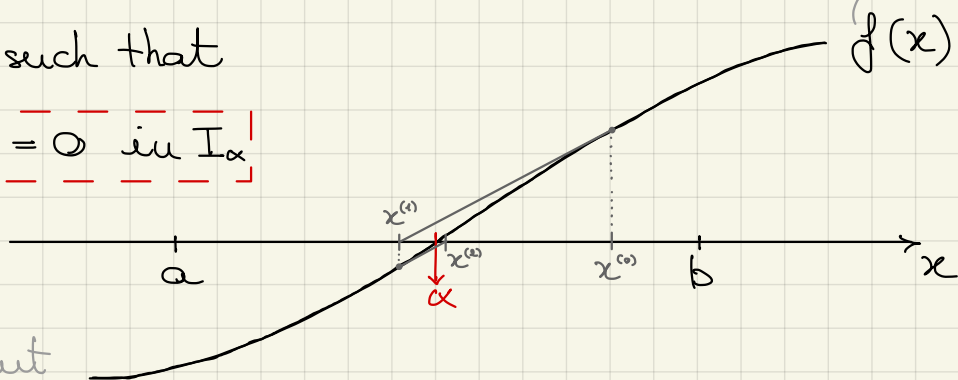
It can be a highly non-linear problem, whose solution (which might not even be unique) must be computed numerically

Newton's method

NON-linear

$I_{\alpha} = [a, b]$  such that

$$\exists! \alpha: f(\alpha) = 0 \text{ in } I_{\alpha}$$



there can be many different  $I_{\alpha}$ , each one having a different solution  $\alpha$

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})} \quad \forall k \geq 0$$

Stopping criterion (or "terminating test"):  $|x^{(k+1)} - x^{(k)}| < \epsilon$

Newton's method is a FIXED-POINT (or "Picard") ITERATION:

$$x^{(k+1)} = T_f(x^{(k)})$$

$k$ : iteration counter  $x^{(k)}$ : iterate

The concept behind fixed-point iterations is to transform the initial problem:

$$\text{find } \alpha \text{ s.t. } f(\alpha) = 0$$

whose solution cannot be obtained analytically, into a new problem:

$$\text{find } \alpha \text{ s.t. } T_f(\alpha) = \alpha$$

whose solution can be obtained through the previously seen iteration.

To guarantee the validity of such transformation it must be:

1.  $f(\alpha) = 0 \iff T_f(\alpha) = \alpha$

2.  $\lim_{k \rightarrow \infty} x^{(k)} = \alpha \iff e^{(k)} = |\alpha - x^{(k)}| \rightarrow \lim_{k \rightarrow \infty} e^{(k)} = 0$

so that for  $x^{(k_0)} = \alpha$  convergence is reached and  $\forall k \geq k_0$   $x^{(k)} = \alpha$  (hence  $\alpha$  is a fixed-point of the iteration  $T_f$ ).

Theorem (existence, uniqueness and convergence of a fixed point of an iteration)

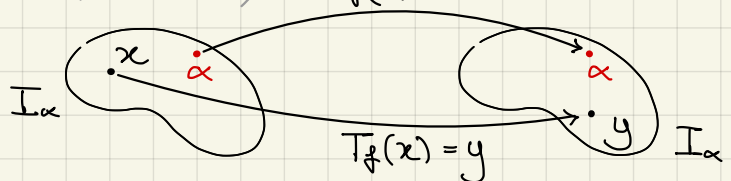
Let  $T_f$  be a fixed-point map (or iteration function) such that:

a)  $T_f: I_\alpha \rightarrow I_\alpha$  (e.g.  $I_\alpha = [a, b] \subset \mathbb{R}$ )

b)  $T_f \in C^0(I_\alpha)$

c)  $T_f \in C^1(I_\alpha)$

d)  $\exists h$  such that  $0 < h < 1$  and  $|T_f'(x)| \leq h \forall x \in I_\alpha$



Then  $\exists! \alpha \in I_\alpha$  such that  $\alpha = T_f(\alpha)$  and

$$\lim_{k \rightarrow \infty} x^{(k)} = \alpha$$

$$\lim_{k \rightarrow \infty} \frac{\alpha - x^{(k+1)}}{\alpha - x^{(k)}} = T_f'(\alpha)$$

The second result of the theorem gives an indication about the speed of convergence of the iteration:

$$|\alpha - x^{(k+1)}| \leq h |\alpha - x^{(k)}| < |\alpha - x^{(k)}| \quad \forall k \geq k_0 > 0$$

i.e.:  $e^{(k+1)} < e^{(k)}$  for  $k$  sufficiently high

$1 > h \geq |T_f'(\alpha)| =$  asymptotic error reduction factor

Newton's method:  $T_f(x) = x - \frac{f(x)}{f'(x)}$   $\rightarrow f \in C^1(I_\alpha), f \in C^2(I_\alpha)$   
 $T_f'(x) = 1 - \frac{(f'(x))^2 - f(x)f''(x)}{(f'(x))^2}$   $f'(\alpha) \neq 0$   
 $= \frac{f(x)f''(x)}{(f'(x))^2}$   
 $T_f'(\alpha) = \frac{f(\alpha)f''(\alpha)}{(f'(\alpha))^2} = 0 \rightarrow$  great reduction factor!  
 $0 =$   $\neq 0$

### Theorem

Assume  $f \in C^2(I_\alpha)$  with  $f'(\alpha) \neq 0$ . Then  $\forall x^{(0)} \in I_\alpha$  the Newton's method converges and it is:

$$\left[ \lim_{k \rightarrow \infty} \frac{\alpha - x^{(k+1)}}{(\alpha - x^{(k)})^2} = C_\alpha = \frac{f''(\alpha)}{2f'(\alpha)} \right]$$

i.e.:  $|e^{(k+1)}| \leq |C_\alpha| |e^{(k)}|^2$  for  $k$  sufficiently high

second order method

Issue of Newton's method: do we always know  $f'(x)$ ?

Typically, just like  $f(x)$ , also  $f'(x)$  cannot be dealt with analytically. Therefore, we would need to somehow approximate  $f'(x)$ , thus introducing inaccuracies in the method.

In conclusion, the high convergence rate of Newton's method comes with the requirement of knowing the derivative of the function for which we want to find the zero. If this information is missing, then the convergence order might be less than the theoretical one.



Back to our problem:

find  $u^*$  such that  $\mathcal{N}(u^*) = 0 \iff u^{(k+1)} = T_x(u^{(k)})$

$$T_x = ?$$

$\mathcal{N}$  is a non-linear differential operator

$$\mathcal{N}: X \rightarrow Y$$

$$u \in X = ?$$

As we will need to use the finite element method, we can expect:

$$X = H^1(\Omega)$$

and also:

$$Y = X = H^1(\Omega)$$

Let's see the implementation of  $T_x$  through Newton's method:

$$u^{(k+1)} = u^{(k)} - \frac{\mathcal{N}(u^{(k)})}{\mathcal{N}'(u^{(k)})}$$

given  $u^{(0)} \in X$

increment  $\leftarrow$  set  $\delta u^{(k)} := u^{(k+1)} - u^{(k)}$

$\implies$   $\left\{ \begin{array}{l} \text{solve: } \mathcal{N}'(u^{(k)}) \delta u^{(k)} = -\mathcal{N}(u^{(k)}) \\ \text{update: } u^{(k+1)} = u^{(k)} + \delta u^{(k)} \end{array} \right.$  residual  $\rightarrow$

$\forall k \geq 0$  until convergence  $\rightarrow |\delta u^{(k)}| \leq \varepsilon$

Note that the expression  $\mathcal{N}'(u^{(k)}) \delta u^{(k)} = -\mathcal{N}(u^{(k)})$  is a linear partial differential equation  $Ax = b$

$\mathcal{N}'(u^{(k)})$  is the so called "Fréchet derivative":

$$\mathcal{N}'(u^{(k)}) \delta u^{(k)} = \varepsilon \frac{\partial I_h(\delta u^{(k)})}{\partial x} + \alpha \delta u^{(k)} + \left[ \frac{\partial x(u^{(k)})}{\partial u} - \frac{\partial g(u^{(k)})}{\partial u} \right] \delta u^{(k)} = -\mathcal{N}(u^{(k)})$$

$$\varepsilon \frac{\partial I_h(\delta u^{(k)})}{\partial x} + \gamma^{(k)} \delta u^{(k)} = -\mathcal{N}(u^{(k)})$$

$$\text{where } \gamma^{(k)} = \alpha + \left[ \frac{\partial x(u^{(k)})}{\partial u} - \frac{\partial g(u^{(k)})}{\partial u} \right]$$

The issue with Newton's method is the necessity to know the derivative of  $x(u)$  and  $g(u)$



Alternatively, we can use another method (i.e. another fixed-point iteration), based on the expression of our problem:

$$\mathcal{N}(u) = \varepsilon \frac{\partial I_{in}(u)}{\partial x} + \frac{u}{\Delta t} - \frac{V_m^n}{\Delta t} + \frac{u}{C_m} g_{tot}(u) - \frac{1}{C_m} \sum_{\alpha} g_{\alpha}(u) E_{\alpha}(u) = 0$$

Idea: compute the new value  $u^{(k+1)}$  using the old value  $u^{(k)}$  for the evaluation of the non-linear terms

given  $u^{(0)} \in X$

$$\text{solve: } \varepsilon \frac{\partial I_{in}(u^{(k+1)})}{\partial x} + \left[ \frac{1}{\Delta t} + \frac{g_{tot}(u^{(k)})}{C_m} \right] u^{(k+1)} = \frac{V_m^n}{\Delta t} + \frac{1}{C_m} \sum_{\alpha} g_{\alpha}(u^{(k)}) E_{\alpha}(u^{(k)})$$

+ boundary conditions ...

$\forall k \geq 0$  until convergence

Defining  $\gamma^{(k)} = \frac{1}{\Delta t} + \frac{g_{tot}(u^{(k)})}{C_m} > 0$  and  $f^{(k)} = \frac{V_m^n}{\Delta t} + \frac{1}{C_m} \sum_{\alpha} g_{\alpha}(u^{(k)}) E_{\alpha}(u^{(k)})$

we get:  $\frac{\partial}{\partial x} \left( -\frac{a}{2g_{\alpha} C_m} \frac{\partial u^{(k+1)}}{\partial x} \right) + \gamma^{(k)} u^{(k+1)} = f^{(k)}$

that resembles a diffusion-reaction problem for which we can express the weak formulation (noting that  $\gamma^{(k)} > 0$ ).

(Remember that  $u^{(k)} = u^{(k)}(x)$  is still a function - a point in the space of functions  $H^1(\Omega)$ )

$$\int_0^L \phi \left[ \frac{1}{2\pi a C_m} \frac{\partial I_{in}(u)}{\partial x} + \gamma^{(k)} u \right] = \int_0^L \phi f^{(k)} \quad \forall \phi \in X = H^1(\Omega)$$

$$(w) \left[ \frac{1}{2\pi a C_m} \int_0^L \phi \vec{I}_{in} \cdot \vec{n} - \frac{1}{2\pi a C_m} \int_0^L I_{in}(u) \frac{\partial \phi}{\partial x} + \int_0^L \gamma^{(k)} u \phi = \int_0^L f^{(k)} \phi \right]$$

$$\text{B.c.: } \vec{I}_{in} \cdot \vec{n} = \alpha u - \beta$$

$$\frac{1}{2\pi a C_m} \left[ \phi(0)(\alpha_0 u(0) - \beta_0) + \phi(L)(\alpha_L u(L) - \beta_L) \right] + \int_0^L \frac{a}{2g_{\alpha} C_m} \frac{\partial u}{\partial x} \frac{\partial \phi}{\partial x} + \int_0^L \gamma^{(k)} u \phi = \int_0^L f^{(k)} \phi$$

$$\left[ \mathcal{B}(u, \phi) = F(\phi) \right] \quad \forall \phi \in X$$

Solvability of the problem and uniqueness of the

solutions are granted like for any other weak formulation, as already seen:

$$B(u, u) = \frac{1}{2\alpha a c_m} [\alpha_0 (u(0))^2 + \alpha_L (u(L))^2 - u(0) \beta_0] + \frac{a}{2 \int_{\min}^{\max} c_m dx} \int_0^L \left( \frac{\partial u}{\partial x} \right)^2 + \int_0^L \gamma^{(k)} u^2 \geq \frac{\|u\|_x^2}{K}$$

$$\frac{1}{K} = \min \left\{ \frac{a}{2 \int_{\max} c_m}, \min_{x \in \Omega} \gamma^{(k)} \right\}$$

$$\frac{\|u\|_x^2}{K} \leq B(u, u) = F(u) = \|f u\|_{L^1} \leq \|f\|_{L^2} \|u\|_{L^2} \leq \|f\|_{L^2} \|u\|_x$$

Cauchy-Schwartz

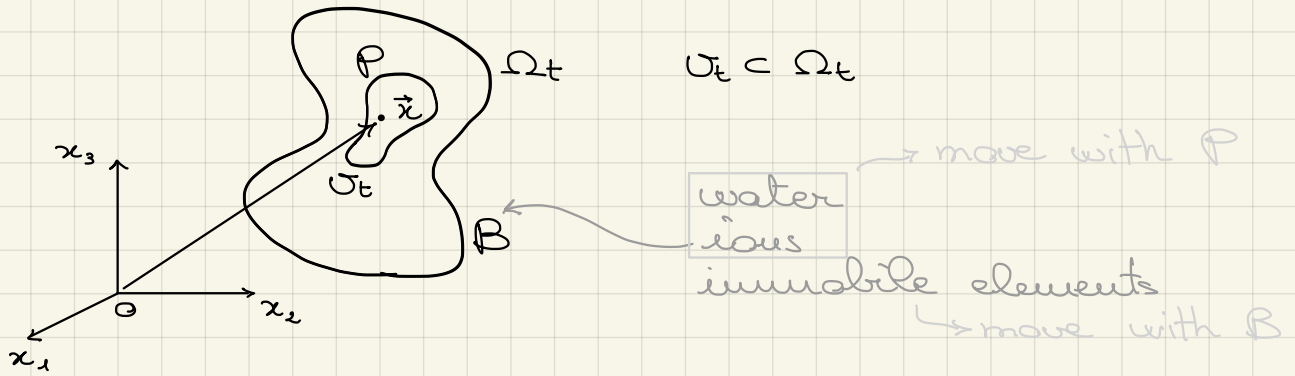
$$\boxed{\|u\|_x \leq K \|f\|_{L^2}}$$

This verification ensures that every step of the iteration embodies a finite element method for the weak formulation of a problem that is well posed.

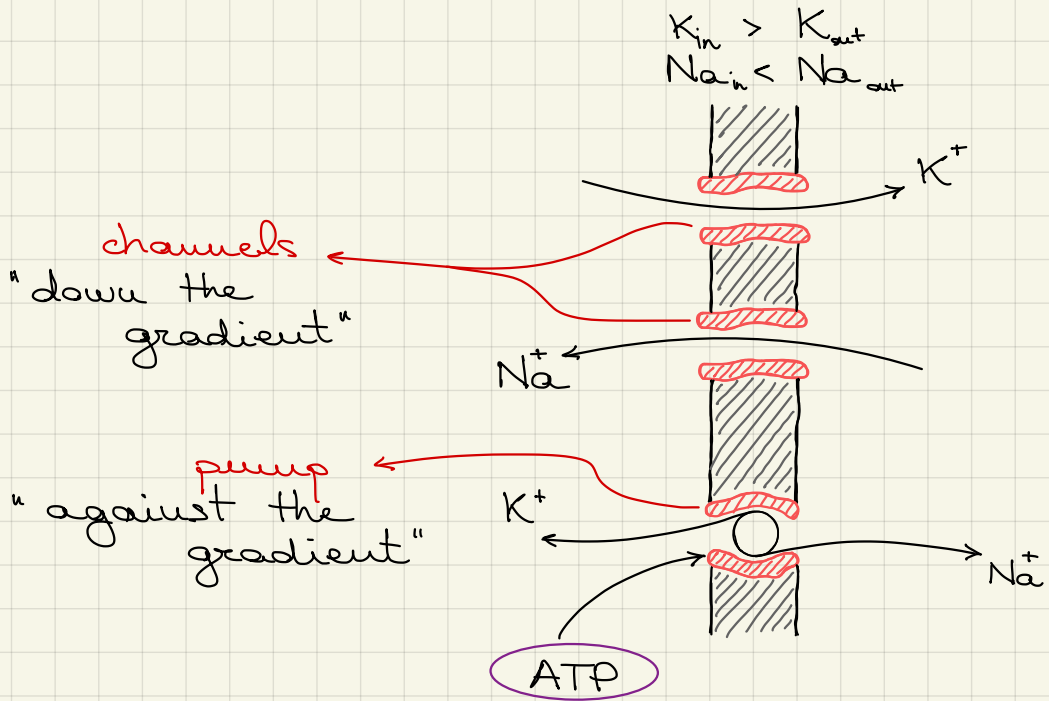
(It does not guarantee, however, that our fixed-point iteration will converge.)

# Velocity-Extended Poisson-Nernst-Planck model

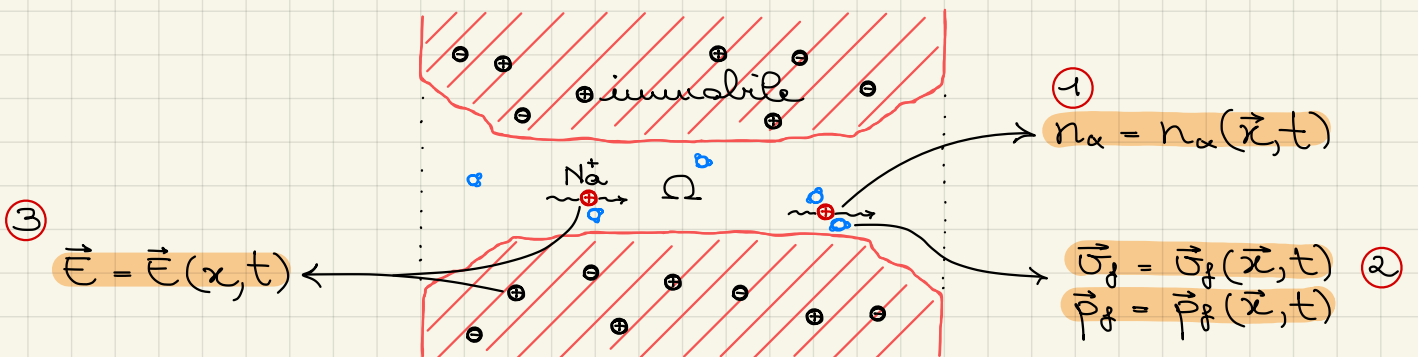
for ion electro-fluid-diffusion



Water and ions are the moving parts of the model that both affect and are affected by whatever variable in play; immobile elements are instead fixed objects whose influence on the system is always the same, since they are in turn not influenced by the environment (like fixed charges in a channel).



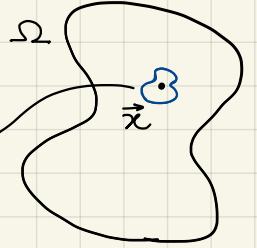
Need to upscale all the microscopic movements of particles into an "average" model: the VE-PNP model satisfies this need.



① What is the number density  $n_\alpha$  of the moving ions, if we look at the microscopic scale?

→ notion of material volume

single charged particle  $dQ \leftrightarrow$  material volume  $dV_t$   
(in point  $\vec{x}$ )



$$dQ_\alpha(\vec{x}, t) = q z_\alpha n_\alpha(\vec{x}, t) dV_t$$

$$\implies n_\alpha(\vec{x}, t) = \frac{1}{q z_\alpha} \frac{dQ(\vec{x}, t)}{dV_t}$$

②  $\vec{v}_f$  and  $\vec{p}_f$  are the velocity vector and pressure field, respectively, associated to the fluid (i.e. water)

③  $\vec{E}$  is the electric field generated by both fixed charges and moving ions themselves.

integral quantity      local quantity

Mass of  $\alpha$  ions in volume  $P$ :  $M_\alpha(P, t) := \int_{V_t} m_\alpha n_\alpha(\vec{x}, t) dV_t$

Mass density of  $\alpha$  ions:  $\rho_\alpha^m(\vec{x}, t) := m_\alpha n_\alpha(\vec{x}, t)$

Charge of  $\alpha$  ions in volume  $P$ :  $Q_\alpha(P, t) := \int_{V_t} q z_\alpha n_\alpha(\vec{x}, t) dV_t$

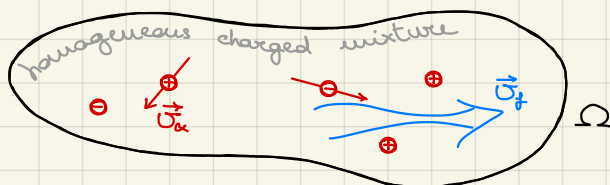
Electric charge density of  $\alpha$  ions:  $\rho_\alpha^e(\vec{x}, t) := q z_\alpha n_\alpha(\vec{x}, t)$

Now that we have all the definitions we need, we have to relate each of them through balance equations.

The domain  $\Omega$  from which we are going to derive these equations is a homogeneous charged mixture (i.e. an ionic solution), whose constituents are in aqueous phase only

→  $M + 1$  constituents ( $M$  ion species + water)  
 $\alpha$  of

Ionic solution



↓  
incompressible, electrically neutral fluid ( $\rho_f^m = \text{const.}$ ,  $\rho_f^e = 0$ )

# Hierarchy of balance equations:

- balance of mass
- " " momentum

focus on these two for the derivation of VE-PNP model

- " " energy → ionic solution is isothermal
- " " angular momentum → automatically satisfied

it is an approx. that might not be good in some real scenarios

## Regarding the ions:

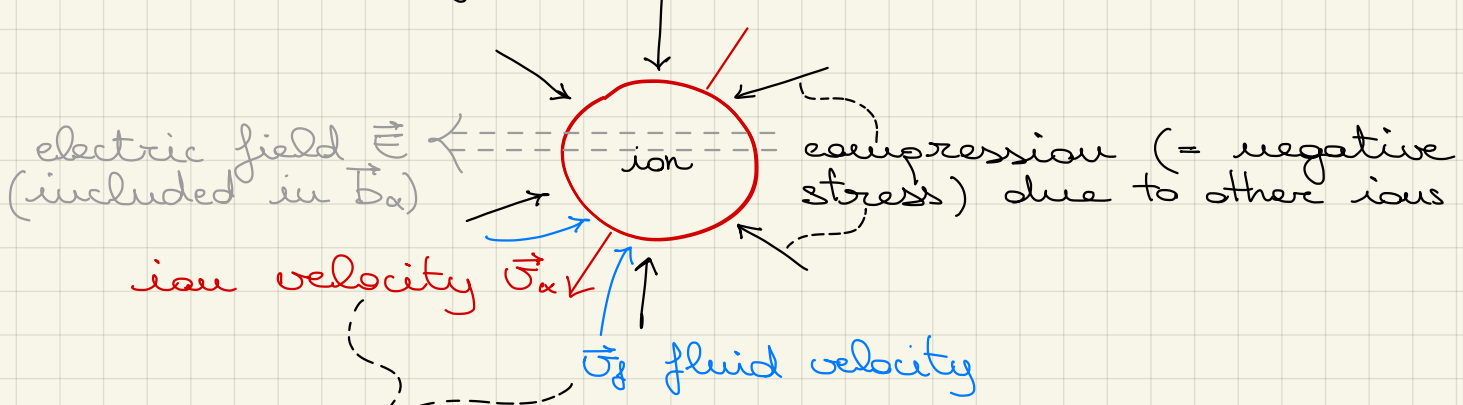
balance of mass:  $\frac{\partial \rho_\alpha^m}{\partial t} = \beta_\alpha - \nabla \cdot (\rho_\alpha^m \vec{v}_\alpha)$   $\left[ \frac{\text{kg}}{\text{m}^3 \text{s}} \right]$

time rate of change of ion mass density → net production rate

balance of forces:  $\frac{\partial (\rho_\alpha^m \vec{v}_\alpha)}{\partial t} + \nabla \cdot (\rho_\alpha^m \vec{v}_\alpha \otimes \vec{v}_\alpha) = \nabla \cdot \underline{T}_\alpha + \vec{b}_\alpha$   $\left[ \frac{\text{N}}{\text{m}^3} \right]$

tensor product  $\underline{a}, \underline{b} \in \mathbb{R}^n: \underline{a} \otimes \underline{b} = [a_i b_j]_{i,j=1}^n \in \mathbb{R}^{n \times n}$  stress tensor volumetric force density

We assume that ions, from the point of view of continuum mechanics, may be regarded as a COMPRESSIBLE fluid:



friction (stress) against fluid and other ions as well

identity matrix ← bulk modulus → shear viscosity

→  $\underline{T}_\alpha = -p_\alpha \underline{I} + \lambda_\alpha \nabla \cdot \vec{v}_\alpha \underline{I} + 2\mu_\alpha \underline{D}(\vec{v}_\alpha)$   $\left[ \frac{\text{N}}{\text{m}^2} \right]$

stress due to pressure (compression) resistance to compression friction

where  $\underline{D}(\vec{v}) = \frac{1}{2}(\underline{J}(\vec{v}) + (\underline{J}(\vec{v}))^T)$  is the "symmetric gradient" and  $\underline{J}(\vec{v})$  is the Jacobian of  $\vec{v}$

↑ temperature

$$p_\alpha = k_B \theta n_\alpha \quad \text{ideal gas law}$$

$$\vec{b}_\alpha = q \vec{E} z_\alpha n_\alpha \quad \text{electric field}$$

$$- \sum_{\substack{\gamma=1 \\ \gamma \neq \alpha}}^M C_{\alpha\gamma} (\vec{v}_\alpha - \vec{v}_\gamma) \quad \text{ion-ion interaction}$$

$$- C_{\alpha f} (\vec{v}_\alpha - \vec{v}_f) \quad \text{ion-fluid interaction}$$

(neglecting gravitational forces)

$$\frac{\partial \rho_\alpha^m}{\partial t} + \vec{\nabla} \cdot (\rho_\alpha^m \vec{v}_\alpha) = \beta_\alpha$$

$$\frac{\partial (\rho_\alpha^m \vec{v}_\alpha)}{\partial t} + \vec{\nabla} \cdot (\rho_\alpha^m \vec{v}_\alpha \otimes \vec{v}_\alpha) = \vec{\nabla} \cdot \underline{\underline{T}}_\alpha + \vec{b}_\alpha$$

Ion model  
 $\alpha = 1 \dots M$

$$\underline{\underline{T}}_\alpha = -k_B \theta n_\alpha \underline{\underline{I}} + \lambda_\alpha \vec{\nabla} \cdot \vec{v}_\alpha \underline{\underline{I}} + 2\mu_\alpha \underline{\underline{D}}(\vec{v}_\alpha)$$

$$\vec{b}_\alpha = q z_\alpha n_\alpha \vec{E} - \sum_{\gamma \neq \alpha} C_{\alpha\gamma} (\vec{v}_\alpha - \vec{v}_\gamma) - C_{\alpha f} (\vec{v}_\alpha - \vec{v}_f)$$

• Regarding the fluid:

balance of mass:  $\vec{\nabla} \cdot \vec{v}_f = 0 \rightarrow$  incompressibility of water ( $\rho_f^m = \text{const.}, \beta_f = 0$ )

balance of forces:  $\rho_f^m \frac{\partial \vec{v}_f}{\partial t} + \rho_f^m \vec{\nabla} \cdot (\vec{v}_f \otimes \vec{v}_f) = \vec{\nabla} \cdot \underline{\underline{T}}_f + \vec{b}_f$

$$\underline{\underline{T}}_f = -p_f \underline{\underline{I}} + 2\mu_f \underline{\underline{D}}_f(\vec{v}_f)$$

no compression mechanisms

$$\vec{b}_f = - \sum_{\alpha=1}^M C_{\alpha f} (\vec{v}_f - \vec{v}_\alpha)$$

neutral fluid

$$\vec{\nabla} \cdot \vec{v}_f = 0$$

$$\rho_f^m \frac{\partial \vec{v}_f}{\partial t} + \rho_f^m \vec{\nabla} \cdot (\vec{v}_f \otimes \vec{v}_f) = \vec{\nabla} \cdot \underline{\underline{T}}_f + \vec{b}_f$$

$$\underline{\underline{T}}_f = -p_f \underline{\underline{I}} + 2\mu_f \underline{\underline{D}}_f(\vec{v}_f)$$

$$\vec{b}_f = - \sum_{\alpha} C_{\alpha f} (\vec{v}_f - \vec{v}_\alpha)$$

Fluid model



- We also need a coupling with electric forces:

Electromagnetic (Maxwell's) equations:

$$\vec{\nabla} \times \vec{E} = - \frac{\partial \vec{B}}{\partial t}$$

$$\vec{\nabla} \times \vec{H} = \vec{j} + \frac{\partial \vec{D}}{\partial t}$$

$$\vec{\nabla} \cdot \vec{D} = \rho^{el}$$

$$\vec{\nabla} \cdot \vec{B} = 0$$

$$\vec{\nabla} \times \vec{E} = 0 \leftrightarrow \vec{E} = -\vec{\nabla} \psi$$

quasi-static approx.  $\rightarrow \vec{H} = 0$

$$\vec{\nabla} \cdot \vec{D} = \rho^{el}$$

where  $\vec{B} = \mu \vec{H}$ ,  $\vec{D} = \epsilon \vec{E}$ .

$$\rho^{el} = \rho_{fixed}^{el}(\vec{x}) + \rho_{mobile}^{el}(\vec{x}, t) = \rho_{fixed}^{el}(\vec{x}) + q \sum_{\alpha=1}^M z_{\alpha} n_{\alpha}(\vec{x}, t)$$

$\epsilon = \epsilon_r \cdot \epsilon_0$  but what  $\epsilon_r$  should we consider?

We will simply assume  $\epsilon = \epsilon_m = \text{const.}$ , where  $\epsilon_m$  gathers the permittivity of all the materials in the system (channel fluid, membrane, etc.) in an effective parameter.

$$\rightarrow \vec{\nabla} \cdot (-\epsilon_m \vec{\nabla} \psi) = \rho_{fixed}^{el} + q \sum_{\alpha=1}^M z_{\alpha} n_{\alpha} \quad \text{Poisson equation}$$

$$\vec{\nabla} \cdot (-\epsilon_m \vec{\nabla} \psi) = \rho_{fixed}^{el} + q \sum_{\alpha} z_{\alpha} n_{\alpha}$$

$$\vec{E} = -\vec{\nabla} \psi$$

Electric model

$I_{ion}$   
&  
Fluid models  
&  
Electric



Fully coupled model  
for a homogeneous  
charged mixture

The VE-PNP model is a simplified form of this fully coupled model:

1)  $\beta_{\alpha} = 0$

2) neglect all inertial terms in the momentum balance equation



3)  $\lambda_\alpha, \mu_\alpha = 0$

4) neglect ion-ion interaction ( $C_{\alpha\beta} = 0$ )

$$\frac{\partial \rho_\alpha^m}{\partial t} + \vec{\nabla} \cdot (\rho_\alpha^m \vec{U}_\alpha) = 0 \quad \textcircled{b}$$

$$\vec{\Theta} = -k_B \theta \vec{\nabla} n_\alpha + q z_\alpha n_\alpha \vec{E} - C_{\alpha f} (\vec{U}_\alpha - \vec{U}_f) \quad \textcircled{a}$$

$$\vec{\nabla} \cdot \vec{U}_f = 0$$

$$\vec{\Theta} = -\vec{\nabla} p_f + 2\mu_f \vec{\nabla} \cdot (\underline{\underline{D}}(\vec{U}_f)) - \sum_\alpha C_{\alpha f} (\vec{U}_\alpha - \vec{U}_f)$$

From  $\textcircled{a}$  we can derive:

$$\vec{U}_\alpha = \vec{U}_f + \frac{1}{C_{\alpha f}} [-k_B \theta \vec{\nabla} n_\alpha + q z_\alpha n_\alpha \vec{E}] \quad \textcircled{1}$$

compare with

$$\vec{J}_\alpha = q z_\alpha n_\alpha \vec{U}_\alpha = q z_\alpha n_\alpha \vec{U}_f + q |z_\alpha| \mu_\alpha^{\text{el}} n_\alpha \vec{E} - q z_\alpha D_\alpha \vec{\nabla} n_\alpha \quad \textcircled{2}$$

which we saw in our previous discussions.

$$\rightarrow C_{\alpha f} = \frac{k_B \theta n_\alpha}{D_\alpha} \quad \text{Stokes' drag theory}$$

$$\rightarrow D_\alpha = \frac{\mu_\alpha^{\text{el}}}{|z_\alpha|} \frac{k_B \theta}{q} \quad \text{Einstein's relation}$$

$$\Rightarrow C_{\alpha f} = \frac{q |z_\alpha| n_\alpha}{\mu_\alpha^{\text{el}}}$$

Multiply  $\textcircled{1}$  by  $q z_\alpha n_\alpha$ , substituting  $C_{\alpha f}$ .

$$\begin{aligned} q z_\alpha n_\alpha \vec{U}_\alpha &= q z_\alpha n_\alpha \vec{U}_f + q z_\alpha n_\alpha \frac{\mu_\alpha^{\text{el}}}{q |z_\alpha| n_\alpha} [-k_B \theta \vec{\nabla} n_\alpha + q z_\alpha n_\alpha \vec{E}] \\ &= q z_\alpha n_\alpha \vec{U}_f - q z_\alpha \frac{\mu_\alpha^{\text{el}}}{|z_\alpha|} \frac{k_B \theta}{q} \vec{\nabla} n_\alpha + \frac{q^2 z_\alpha^2 n_\alpha \vec{E}}{q |z_\alpha|} \\ &= q z_\alpha n_\alpha \vec{U}_f - q z_\alpha D_\alpha \vec{\nabla} n_\alpha + q |z_\alpha| n_\alpha \vec{E} \end{aligned}$$

which is exactly equation  $\textcircled{2}$

Now look at equation  $\textcircled{b}$ . Since  $\rho_\alpha^m = m_\alpha n_\alpha$  it is:

$$\frac{\partial n_\alpha}{\partial t} + \vec{\nabla} \cdot (n_\alpha \vec{U}_\alpha) = 0$$

Multiplying both sides by  $qz_\alpha$  we get:

$$qz_\alpha \frac{\partial n_\alpha}{\partial t} + \vec{\nabla} \cdot \vec{j}_\alpha = 0$$

So the final form of VE-PNP model is

$$qz_\alpha \frac{\partial n_\alpha}{\partial t} + \vec{\nabla} \cdot \vec{j}_\alpha = 0$$

$$\vec{j}_\alpha = qz_\alpha n_\alpha \vec{v}_\alpha + q|z_\alpha| \mu_\alpha^\text{el} n_\alpha \vec{E} - qz_\alpha D_\alpha \vec{\nabla} n_\alpha$$

$$\vec{\nabla} \cdot \vec{v}_f = 0$$

$$-\vec{\nabla} p_f + 2\mu_f \vec{\nabla} \cdot (\underline{\underline{D}}(\vec{v}_f)) - \sum_\alpha C_{\alpha f} (\vec{v}_\alpha - \vec{v}_f) = \vec{0}$$

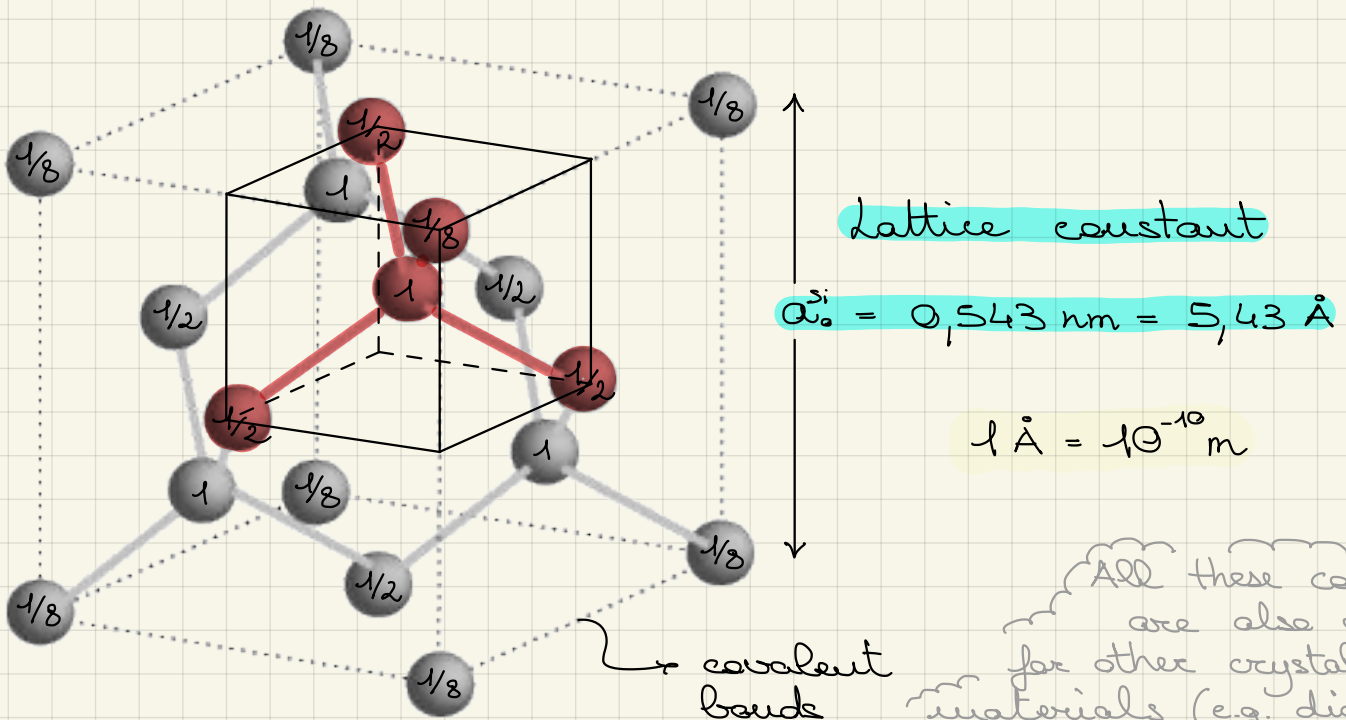
+

$$\vec{\nabla} (-\epsilon_m \vec{\nabla} \psi) = \rho_\text{fixed} + q \sum_\alpha z_\alpha n_\alpha$$

$$\vec{E} = -\vec{\nabla} \psi$$

# Microelectronics

It is of paramount importance understanding the physical scale at which electronic devices operate.



## Silicon - crystal lattice of a unit cell

All these concepts are also valid for other crystalline materials (e.g. diamond, germanium) taking into account different values of  $a_0$ .

Solid silicon has a periodic structure (crystal) whose primary component is a tetrahedron of four atoms. These tetrahedra bond together to form the unit cell, which is the periodic element (i.e. the one that is repeated identically) of the crystal.

Atom density:  $\frac{\# \text{ atoms in a unit cell}}{\text{unit cell volume}}$

$$\# \text{ atoms in a unit cell} = 4 + 6 \cdot \frac{1}{2} + 8 \cdot \frac{1}{8} = 8$$

↓ within the cell    
 ↓ on the faces    
 ↓ at the corners

unit cell volume =  $a_0^3$

⇒ Silicon atom density:  $n_a^{\text{Si}} = 5 \cdot 10^{22} \text{ cm}^{-3}$

(Distance between nearest atoms:  $\frac{\sqrt{3} a_0}{4} = 2,35 \text{ \AA}$ )

Suppose we have an electron travelling through the reticle.

Can we still see it as a particle, even in a sub-nanometer scale ( $a_0 < 1\text{nm}$ )?



Wave-particle duality

- Particle view: Newton's (macroscopic) law of motion

$$\vec{F} = m_e \cdot \vec{a}$$

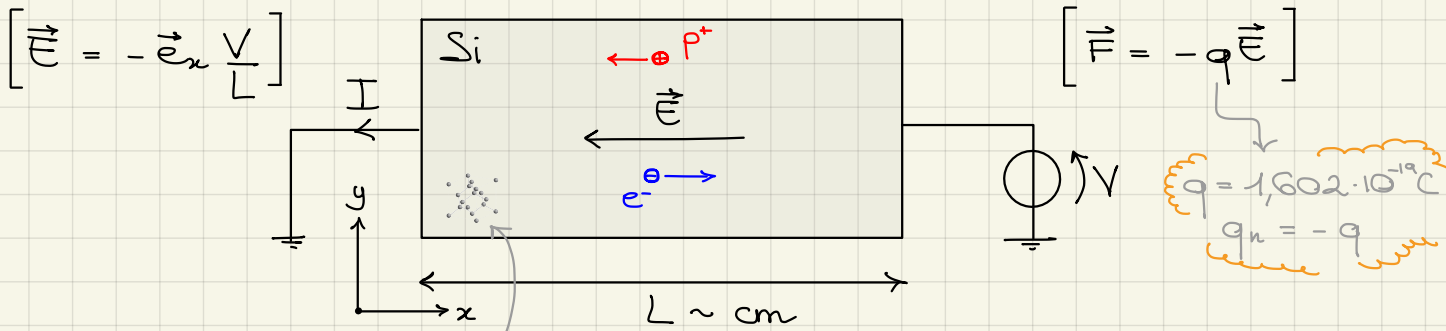
subscript 'n' and 'p' refer to negative and positive charged particles

$m_e = m_0 \cdot \gamma = m_n^*$  effective electron mass

$m_0$ : electron rest mass

might differ depending on the environment

What is typically the source of  $\vec{F}$ ?



yet micro scale is still involved!

from micro to macro scale

$$[-q\vec{E} = m_n^* \vec{a}]$$

Need to merge macroscopic and microscopic phenomena

→ the electron collides (as if it was solid body) with the lattice atoms, whose mass  $m_a$  is much larger than that of the electron. This interaction takes the name of SCATTERING.

Supposing that the electron, accelerated by the electric field  $\vec{E}$ , collides with an atom on average every  $\tau$  losing its kinetic energy in the process, we can

assume that on average it will move with a constant velocity  $\vec{v}_D$ , since some times it is accelerating (no collision) and some times it is at rest (collision).

We can then approximate the acceleration with the ratio of these two quantities:

it's an "effective" acceleration  $\left[ a \sim \frac{v_D}{\tau} \right]$  collision time approx.

$$\implies -q\vec{E} = m_n^* \frac{\vec{v}_D}{\tau} \implies \left[ \vec{v}_D = - \frac{q\tau}{m_n^*} \vec{E} \right]$$

electron mobility  $\mu_n^{el}$

In general, we should consider an ensemble of electrons rather than just one.

electron density  $n$

Now we can compute the current density (associated to electrons):

$$\vec{J}_n = -qn\vec{v}_D = -qn(-\mu_n^{el}\vec{E}) = \underbrace{qn\mu_n^{el}}_{\text{macroscopic}} \vec{E}$$

electric conductivity  $\sigma_n$  microscopic

Remember that for each electron (in un-doped material) there is also a hole, whose equivalent charge is  $+q$ , that participates to the overall current density.

Hence we can define a  $\mu_p^{el}$  and  $\sigma_p$  associated to holes, as well as a current density  $\vec{J}_p$  of holes only.

intrinsic concentration of conductive carriers  $\leftarrow n_i^{Si} = 8,2 \cdot 10^9 \text{ cm}^{-3} = n = p$

$$\mu_n^{Si} = 0,14 \frac{\text{m}^2}{\text{Vs}}$$

$$\mu_p^{Si} = 0,048 \frac{\text{m}^2}{\text{Vs}}$$

$$\sigma^{Si} = 2,5 \cdot 10^{-4} \frac{\text{A}}{\Omega\text{m}} = qn_i\mu_n^{Si} + qn_i\mu_p^{Si}$$

(at room temperature)

• Wave view: Schrödinger equation...

Why, in the first place, do we need an alternate description of the "identity" of a travelling electron?

Because Newton's mechanics does not "get along" with the measurement of microscopical phenomena. The sole particle description cannot explain in its entirety the behaviour of, for example, a travelling electron when you try to measure its velocity or position.

Need of travelling wave interpretation:

$$u(z, t) = \sum_r a_r e^{-i(\omega_r t - k_r z)} \quad \text{"train of waves"}$$

$\downarrow$  frequency       $\downarrow$  wave number

For  $\gamma$  large enough:

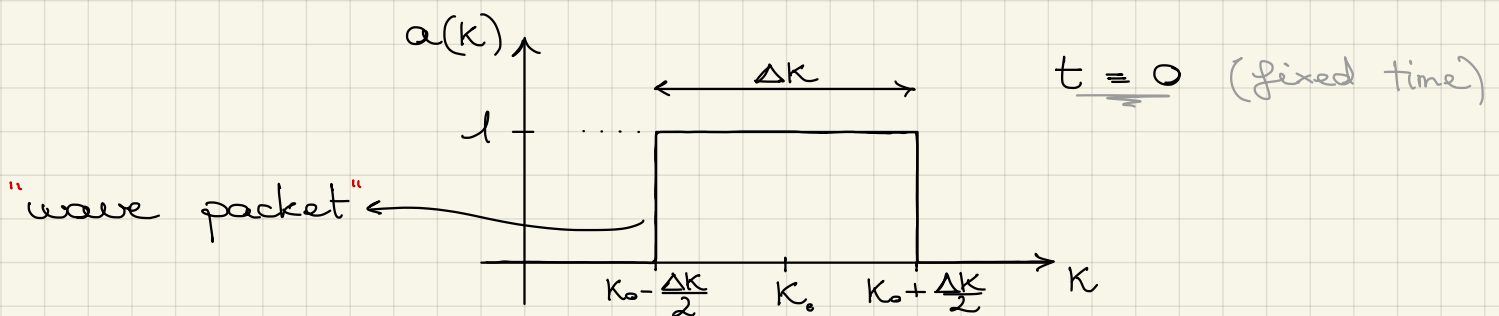
$$u(z, t) = \int_{-\infty}^{+\infty} a e^{-i(\omega t - kz)} dk \quad \left\{ \omega = 2\pi f, \quad k = \frac{2\pi}{\lambda}, \quad \omega = ck \right\}$$

where both  $a$  and  $\omega$  are functions of  $k$ .

$a = a(k) \longrightarrow$  "amplitude dispersion"

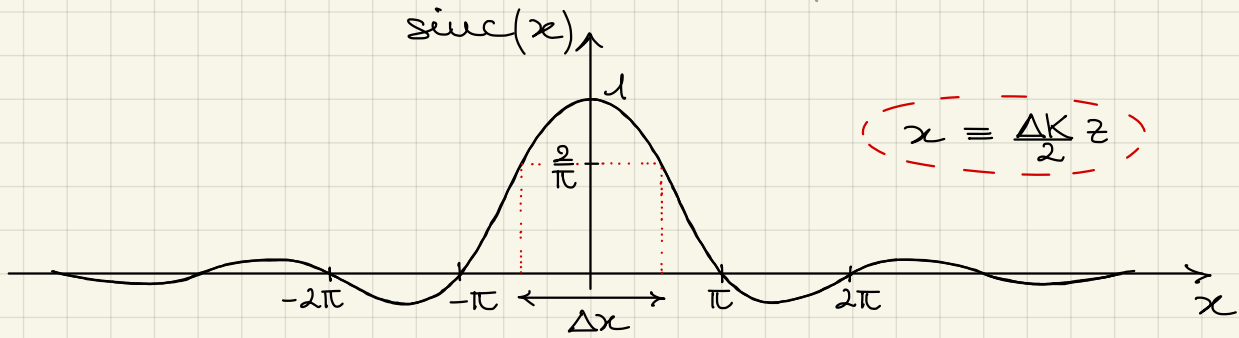
$\omega = \omega(k) \longrightarrow$  "frequency dispersion"  
dispersion equation

These travelling waves (along the  $z$  axis) are oscillating. To show this property, consider an amplitude dispersion as the one below:



$$\begin{aligned}
 u(z, 0) &= \int_{k_0 - \frac{\Delta k}{2}}^{k_0 + \frac{\Delta k}{2}} e^{+ikz} dk = \frac{1}{iz} \left[ e^{ikz} \right]_{k_0 - \frac{\Delta k}{2}}^{k_0 + \frac{\Delta k}{2}} = \frac{e^{i(k_0 + \frac{\Delta k}{2})z} - e^{i(k_0 - \frac{\Delta k}{2})z}}{iz} \\
 &= \frac{\Delta k}{\Delta k \cdot z} \frac{e^{i k_0 z} (e^{i \frac{\Delta k}{2} z} - e^{-i \frac{\Delta k}{2} z})}{2i} = \Delta k e^{i k_0 z} \frac{\sin(\frac{\Delta k}{2} z)}{\frac{\Delta k}{2} z}
 \end{aligned}$$

$$\rightarrow u(z) = \underbrace{\Delta K e^{ik_0 z}}_{\text{travelling wave}} \underbrace{\text{sinc}\left(\frac{\Delta K}{2} z\right)}_{\text{modulating envelop}}$$



Def (amplitude of a wave packet):

it is determined by the point on the  $x$ -axis at which  $\text{sinc}(x)$  changes from 1 to  $\frac{2}{\pi} = \frac{1}{\pi/2} \approx 0,63$

$$\Delta x = \frac{\Delta K}{2} \cdot \Delta z \quad \text{amplitude of the wave packet}$$

$$\text{By definition: } \text{sinc}\left(\frac{\Delta x}{2}\right) = \frac{2}{\pi} \rightarrow \boxed{\frac{\Delta x}{2} = \frac{\pi}{2}}$$

$$\frac{\Delta K}{2} \Delta z = \pi$$

$$\boxed{\Delta K \Delta z = 2\pi}$$

If we want to measure the position of the wave packet with high precision (small  $\Delta z$ ) then the amplitude dispersion grows (big  $\Delta K$ )

Heisenberg indetermination principle

wave packet  $\longleftrightarrow$  electron

In Classical Mechanics:  $p = m \cdot v$   
momentum

In Quantum Mechanics:  $p = \hbar \cdot K$  where  $\hbar = \frac{h}{2\pi}$

$$\Delta p = \hbar \Delta K = \frac{h}{2\pi} \Delta K$$

$$\boxed{\Delta p \Delta z = h}$$

reduced Planck's constant

If we want to measure the position of the electron with high precision (small  $\Delta z$ ) then quantifying its momentum becomes harder (big  $\Delta p$ )



## Schrödinger equation

$$i\hbar \frac{\partial \Psi}{\partial t} = -\frac{\hbar^2}{2m} \Delta \Psi + V \Psi$$

Laplacian operator

known  $\Psi = \Psi(\vec{x}, t)$   $V = V(\vec{x}, t)$  given

$$\vec{x} \in \Omega \subseteq \mathbb{R}^d, \quad d \geq 1$$

$\Psi$  is the wave function that gives the probability of finding a particle in position  $\vec{x}$ , at a certain time  $t$ .

Being a probability function,  $\Psi$  satisfies the normalization condition:

$$\int_{\Omega} |\Psi(\vec{x}, t)|^2 d\Omega = 1$$

$V$  is a given potential (energy) field, typically obtained from the solution of a Poisson equation in the same domain  $\Omega$ .

### Example 0:

Assume now it is possible to apply separation of variables to  $\Psi$ :

$$\Psi(\vec{x}, t) = \psi(\vec{x}) \omega(t)$$

Then we can re-write Schrödinger equation as:

$$\psi(\vec{x}) i\hbar \frac{d\omega(t)}{dt} = \omega(t) \left[ -\frac{\hbar^2}{2m} \Delta \psi(\vec{x}) + V(\vec{x}, t) \psi(\vec{x}) \right]$$

Also assume that  $V = V(\vec{x})$  does not vary with  $t$ .

$$\underbrace{\frac{i\hbar}{\omega(t)} \frac{d\omega(t)}{dt}}_{f(t)} = \underbrace{-\frac{\hbar^2}{2m} \frac{\Delta \psi(\vec{x})}{\psi(\vec{x}) + V(\vec{x})}}_{g(\vec{x})} = \underbrace{E}_{\substack{\text{energy of} \\ \text{the particle}}}$$

The only way for  $f(t)$  and  $g(\vec{x})$  to be equal is to be both equal to the same constant  $E$ .

$$\left\{ \begin{array}{l} i\hbar \frac{d\omega(t)}{dt} = E \omega(t) \rightarrow \omega(t) = \bar{\omega} e^{-iEt/\hbar} \\ -\frac{\hbar^2}{2m} \Delta \psi(\vec{x}) + V(\vec{x}) \psi(\vec{x}) = E \psi(\vec{x}) \end{array} \right.$$

eigenvalue problem in differential form  
eigenvalue eigenvector

An eigenvalue-eigenvector problem in differential form is extremely hard to solve and is very costly from a computational standpoint (for a finer discretization of the problem, the computational effort becomes enormously larger).

This shows how hard it is to solve Schrödinger equation, even after all the simplifications we made.

### Example 1: Free Electron Model

1D domain ( $\Omega \subseteq \mathbb{R}$ ,  $d=1$ ) for simplicity.

Electrons in lattice of a metal conductor:

$$V(x) = 0$$

since potential within a conductor is constant e.g. equal to zero. In other words, valence electrons (those occupying the outer orbital in the atoms) are free to move within the material in the so called "sea of electrons".

Under these assumptions, on top of previous simplifications, Schrödinger equation becomes:

$$-\frac{\hbar^2}{2m} \psi''(x) = E \psi(x)$$

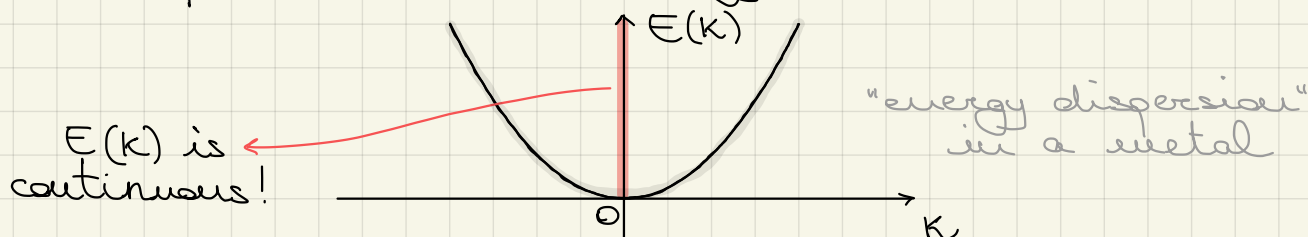
Introducing  $k^2 := \frac{2mE}{\hbar^2}$  we write:

$$\psi''(x) + k^2 \psi(x) = 0$$

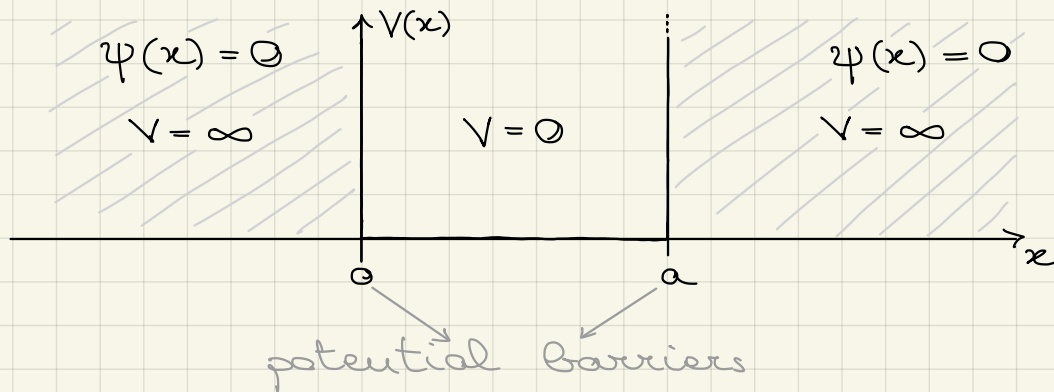
Note that  $k$  corresponds to the wave number. As a matter of fact, the following relation holds:

$$E = \frac{\hbar^2 k^2}{2m} = \frac{|p|^2}{2m}$$

which is indeed the expression for kinetic energy (in fact, no potential energy is present being the potential zero, hence the electron energy corresponds to kinetic energy alone)



## Example 2: Electron Confinement



Electrons are confined in a limited region by two potential barriers of infinite height (even a barrier as high as  $5\text{eV}$  can be considered as infinite since the probability of an electron crossing it  $\sim e^{-\frac{5\text{eV}}{k_B T}} \approx 0$ ). Hence  $\psi(x) = 0$  outside  $(0, a)$ .

The current problem is the same as before ( $V=0$ ) with the addition of boundary conditions in  $x=0$  and  $x=a$ :

$$\psi'' + k^2 \psi = 0$$

$$P(\lambda) = \lambda^2 + k^2 = 0$$

$$\lambda_{1,2} = \pm i k$$

$$\rightarrow \psi(x) = A e^{ikx} + B e^{-ikx}$$

$$\text{b.c.: (1) } \psi(0) = 0$$

$$A = -B$$

$$\begin{aligned} \rightarrow \psi(x) &= A (e^{ikx} - e^{-ikx}) \\ &= 2i \sin(kx) \end{aligned}$$

$$(2) \psi(a) = 0$$

$$\sin(ka) = 0$$

$$\left| k = n \frac{\pi}{a} \right| \quad n = 1, 2, \dots$$

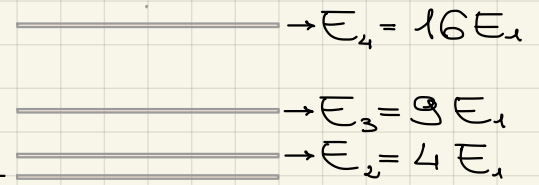
$$\rightarrow \psi(x) = 2i \sin\left(\frac{n\pi}{a} x\right)$$

Note that now the energy distribution is NOT continuous, since  $k^2$  is quantized!

$$E = E_n = \frac{n^2 \hbar^2 \pi^2}{2ma^2}$$

ground level

$$\frac{\hbar^2 \pi^2}{2ma^2} = E_1$$

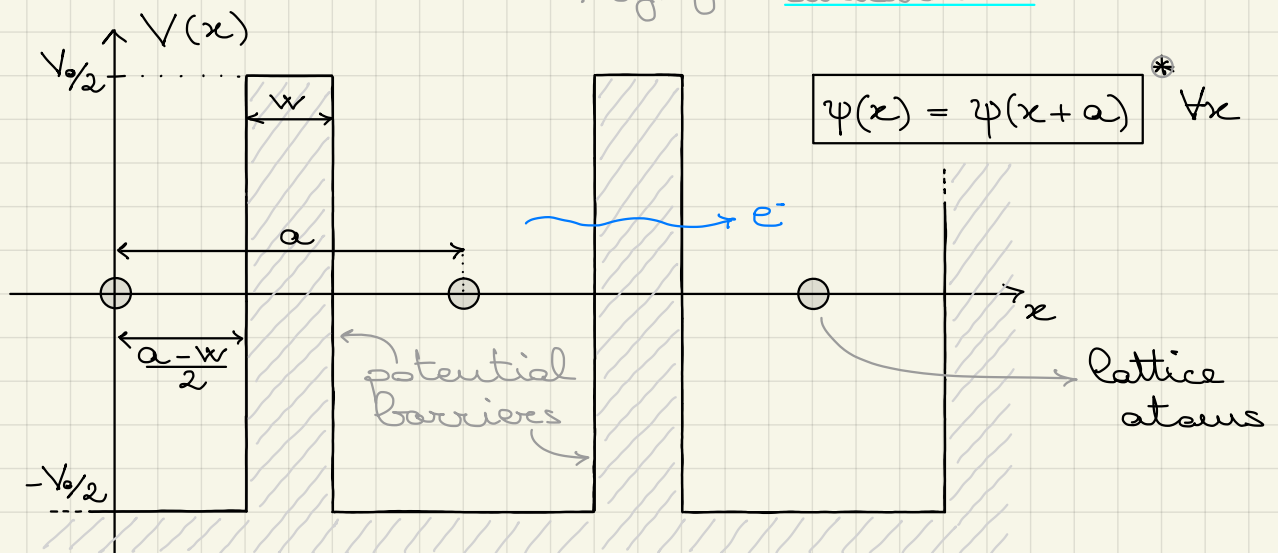


### Example 3: Kronig-Penney model

for the Energy Band theory

It is a more accurate model of the electron confinement example, where energy barriers occur for every atom in the lattice and with finite height.

↳ e.g. of a semiconductor



It can be demonstrated that periodicity\* of the wave function can be achieved if and only if the following condition holds:

$$\cos(Ka) = P \frac{\sin(\alpha a)}{\alpha a} + \cos(\alpha a)$$

where  $P = \frac{ma}{\hbar^2} V_0 w = [\text{pure number}]$

$$\alpha = \sqrt{\frac{2mE}{\hbar^2}} = [m^{-1}]$$

Note: 1) For  $V_0 = 0$  it is  $K = \alpha$

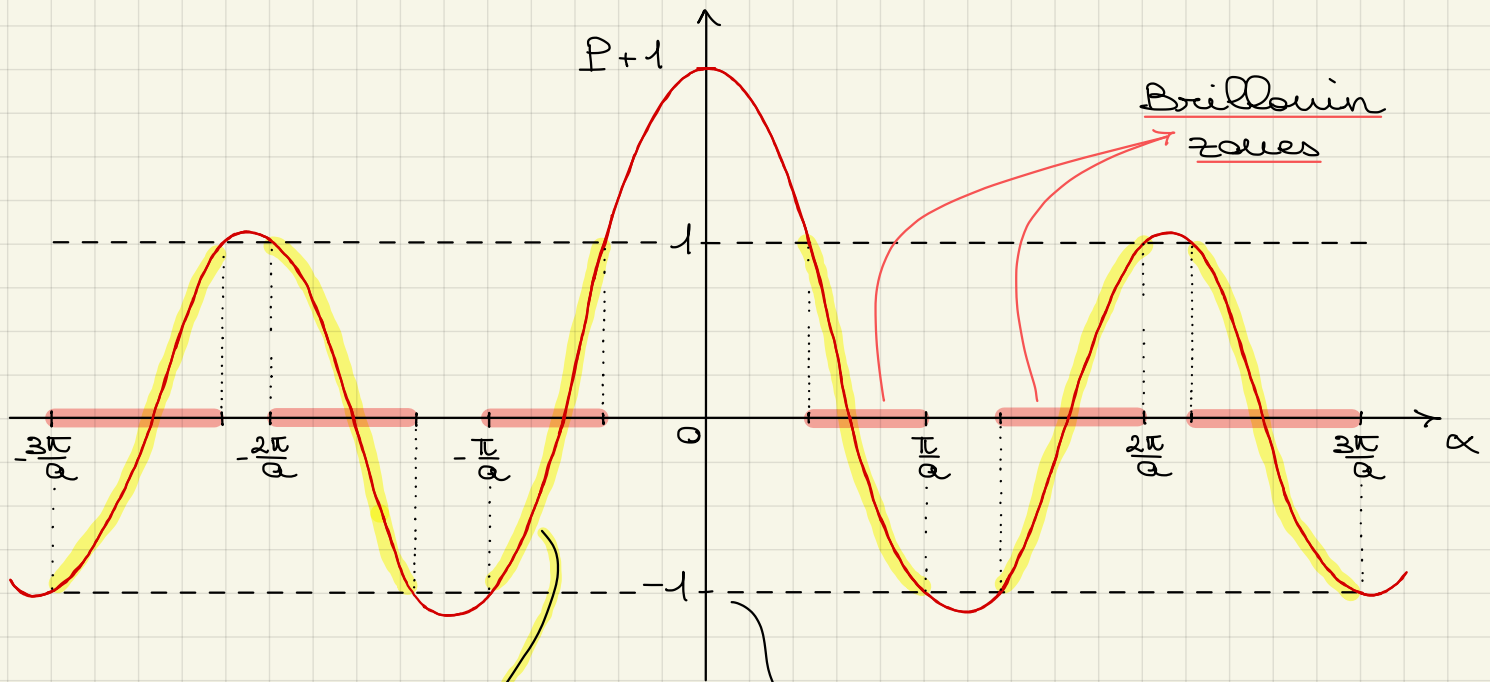
→ Free electron

2) For  $V_0 = \infty$  it is  $0 = \sin(\alpha a)$

→ Electron confinement

Let us plot this non-linear expression of  $k$ :

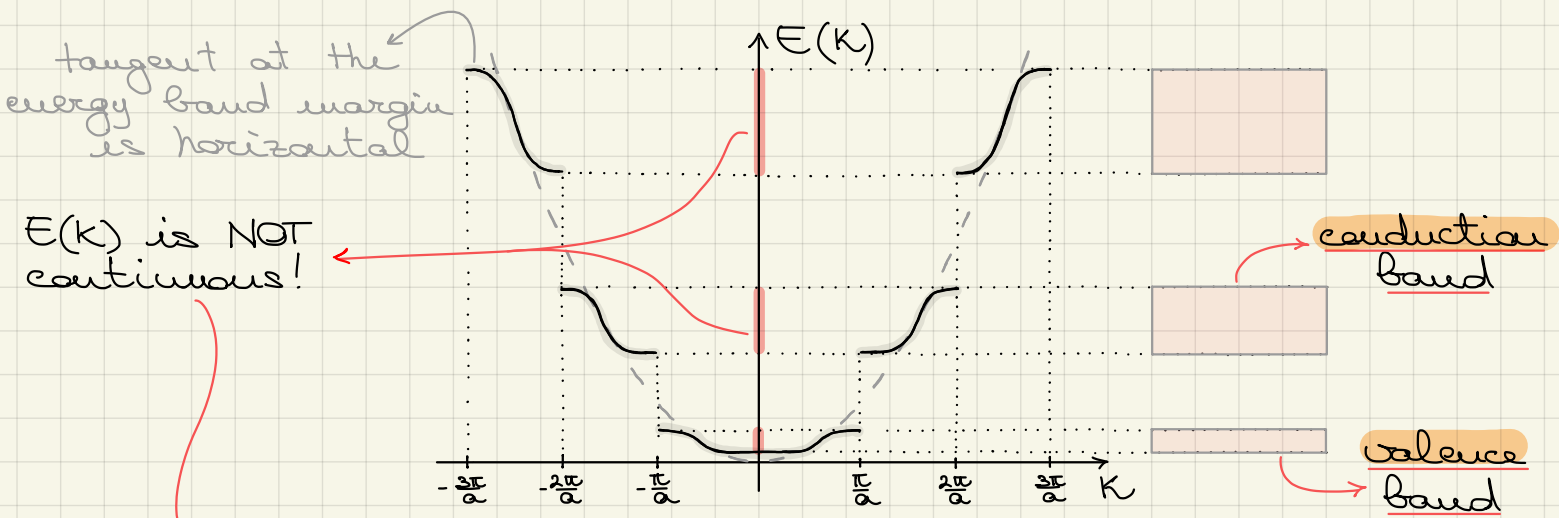
$$k = \frac{1}{a} \arccos \left[ P \frac{\sin(\alpha a)}{\alpha a} + \cos(\alpha a) \right]$$



$a \cdot k$  is equal to the arccos of these lines, only for specific values of  $\alpha$

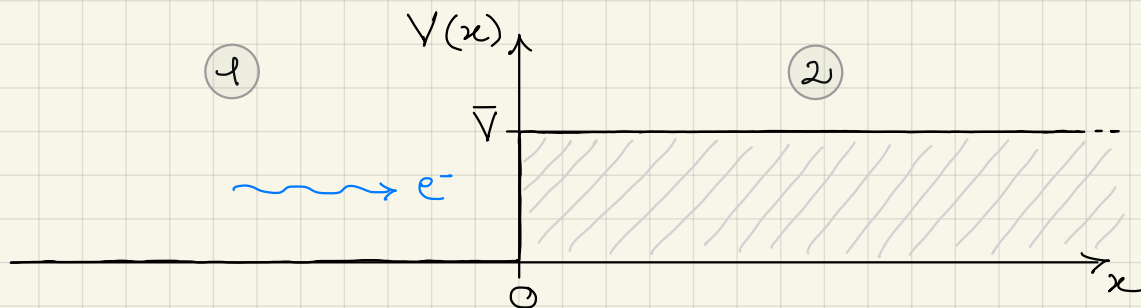
for  $k$  to exist, the argument of the arccos must be between  $-1$  and  $1$

Keeping in mind the expression of  $\alpha$  as a function of  $E$ , it is possible to plot the energy dispersion graph for the Kronig-Penney model:



energy bands where energy is allowed  
energy gaps where energy is forbidden

Example 4: electron collision with finite potential step



$$\psi(x) = \begin{cases} \psi_1(x), & x < 0 \quad (V \equiv 0) \\ \psi_2(x), & x > 0 \quad (V \equiv \bar{V}) \end{cases}$$

①  $\psi_1''(x) + \frac{2mE}{\hbar^2} \psi_1(x) = 0$

$\rightarrow \psi_1(x) = A e^{ikx} + B e^{-ikx}$  with  $k^2 := \frac{2mE}{\hbar^2}$   
 incident wave  $\leftarrow$  reflected wave

②  $\psi_2''(x) + \frac{2m(E - \bar{V})}{\hbar^2} \psi_2(x) = 0$

Two cases. **a**  $E < \bar{V}$  and **b**  $E > \bar{V}$

**a** Let  $\alpha^2 := \frac{2m}{\hbar^2} (\bar{V} - E) > 0$ . Then we can write:

elliptic equation  $\rightarrow \psi_2''(x) - \alpha^2 \psi_2(x) = 0$   
 $\rightarrow \psi_2(x) = C e^{-\alpha x} + D e^{\alpha x}$

To find the exact solution  $\psi(x)$ , impose:

- normalization condition:  $\int_{-\infty}^{+\infty} |\psi(x)|^2 dx = 1$

- transmission conditions:  $\psi(0^-) = \psi(0^+)$   
 $\psi'(0^-) = \psi'(0^+)$

$$\Rightarrow \begin{cases} D = 0 \text{ (otherwise } \psi_2 \text{ would} \\ \text{explode and normalization would fail)} \\ A + B = C \\ Aik - Bik = -\alpha C \end{cases}$$

$$\Rightarrow B = \frac{ik + \alpha}{ik - \alpha} A, \quad C = \frac{2ik}{ik - \alpha} A$$



$$\Rightarrow \psi(x) = \begin{cases} A \left( e^{ikx} + \frac{ik+\alpha}{ik-\alpha} e^{-ikx} \right) & x < 0 \\ A \frac{2ik}{ik-\alpha} e^{-\alpha x} & x > 0 \end{cases}$$

reflection coeff.  $\Gamma$

transmission coeff.  $T$

$$\lim_{\substack{\bar{V} \rightarrow +\infty \\ \alpha \rightarrow +\infty}} \psi(x) = \begin{cases} A (e^{ikx} - e^{-ikx}) = 2i A \sin(kx) & x < 0 \\ 0 & x > 0 \end{cases}$$

The result is interesting since in CM if the potential barrier is higher than the energy of the electron ( $\bar{V} > E$ ) then there is no possibility of ending up in region (2). QM shows instead that there is a small probability of finding the electron beyond the barrier; this probability quickly vanishes the higher is the potential step ( $\bar{V} \rightarrow +\infty$ ) or the farther is the electron from the interface ( $x \rightarrow +\infty$ ).

(b) Let  $(k')^2 := \frac{2m}{\hbar^2} (E - \bar{V})' > 0$ . Then we can write:

$$\psi_2''(x) + (k')^2 \psi_2(x) = 0$$

$$\Rightarrow \psi_2(x) = A' e^{ik'x} + B' e^{-ik'x}$$

Applying same conditions of previous case:

$$\Rightarrow \begin{cases} B' = 0 & \text{(since there is nothing that can produce a reflected wave beyond the barrier)} \\ A + B = A' \\ Aik - Bik = ik'A' \end{cases}$$

$$\Rightarrow B = \frac{k-k'}{k+k'} A, \quad A' = \frac{2k}{k+k'} A$$

$$\Rightarrow \psi(x) = \begin{cases} A \left( e^{ikx} + \frac{k-k'}{k+k'} e^{-ikx} \right) & x < 0 \\ A \frac{2k}{k+k'} e^{ikx} & x > 0 \end{cases}$$

$\Gamma$

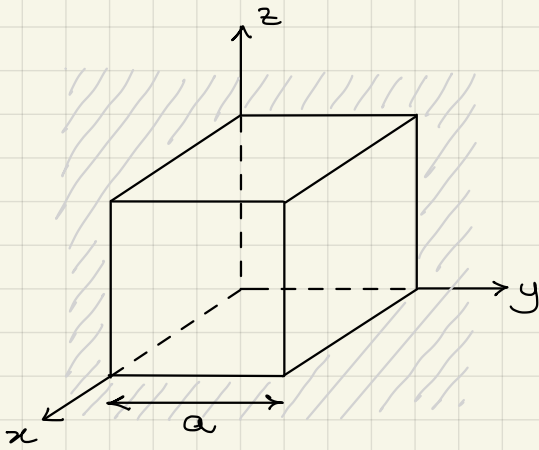
$T$

$$\lim_{\substack{\bar{V} \rightarrow 0 \\ k' \rightarrow k}} \psi(x) = A e^{ikx} \quad \forall x$$

Also this result is interesting since in CM if the potential barrier is lower than the energy of the electron ( $\bar{V} < E$ ) there is no possibility of being reflected. QM shows instead that a reflected electron can still be found in region (1), with a probability that decreases the lower is the potential step ( $\bar{V} \rightarrow 0$ ).



# States density in metals



## 3D electron confinement

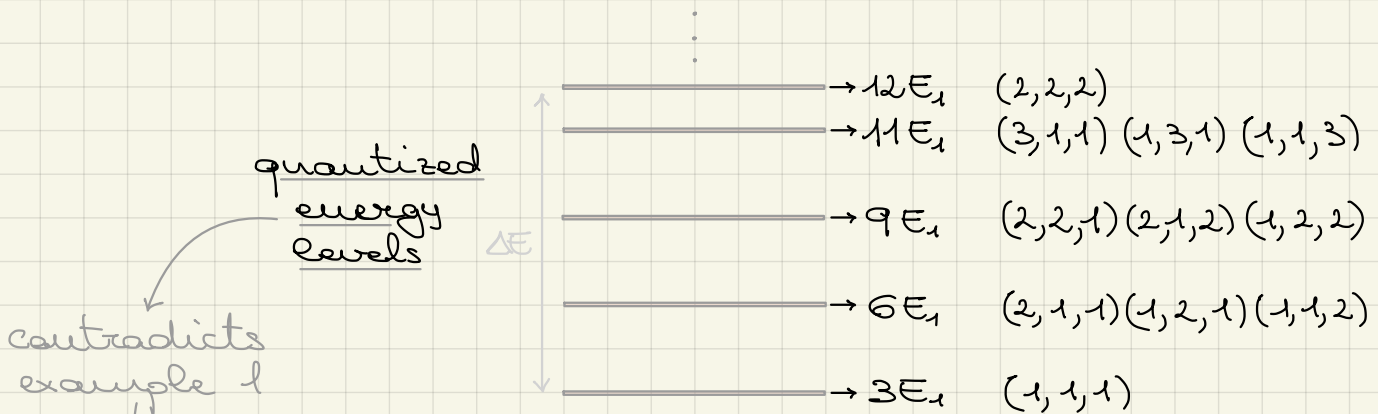
Electrons are forced to stay within a box of volume  $a^3$ . (it represents electrons in a cubic shaped metal).

In analogy with the 1D example:

$$E = \frac{\hbar^2}{2m} (k_x^2 + k_y^2 + k_z^2) = \frac{\pi^2 \hbar^2}{2ma^2} (n_x^2 + n_y^2 + n_z^2)$$

$$\begin{aligned} n_x &= 1, 2, \dots \\ n_y &= 1, 2, \dots \\ n_z &= 1, 2, \dots \end{aligned}$$

$(n_x, n_y, n_z)$  are the so called "Miller indices".



Even if energy is not continuous, for large values of  $E$ , i.e. when the confinement box is much bigger than the atomic scale, the difference between consecutive levels becomes relatively so small that it is then appropriate to speak of a density of energy states as if energy was continuous.

E.g.:  $a = 1\text{cm} = 10^{-2}\text{m} \rightarrow E_1 = \frac{\pi^2 \hbar^2}{2m_n^* a^2} \approx 1,5 \cdot 10^{-14}\text{eV}$

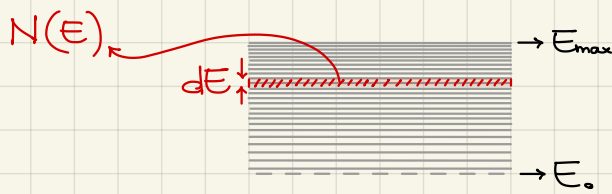
$E_m = 3E_1 m^2$  ( $n_x = n_y = n_z = m$ ) above this level, electrons will just fly out of the box

it's actually assuming  $E_{\text{max}} = 1\text{eV} \rightarrow m_{\text{max}} \sim 5 \cdot 10^6$

more than the difference between consecutive levels

$\Delta E = E_{m+1} - E_m = 3E_1 (2m+1) \rightarrow \Delta E_{\text{max}} \sim 4,5 \cdot 10^{-7}\text{eV} \ll E_{\text{max}}$

hence it is licit to see such highly dense discrete energy distribution as a continuum

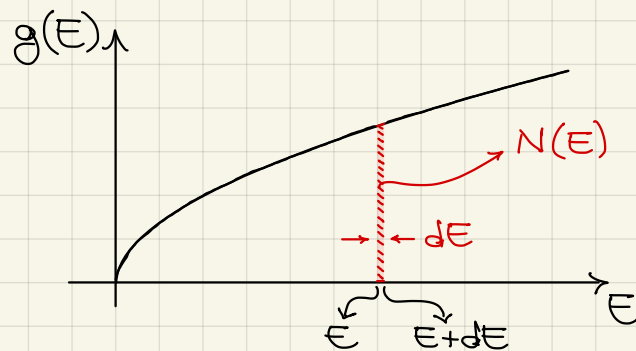


$N(E)$  is the number of energy states within the interval  $(E, E+dE)$  per unit volume

It can be demonstrated that  $N(E)$  is given by the following relation:

$$[N(E) = g(E) dE]$$

where  $[g(E) = \frac{4\pi\sqrt{2Em^3}}{h^3}]$  is the energy states density per unit volume



We now have the number of available states that an electron is allowed to occupy in our metal box.

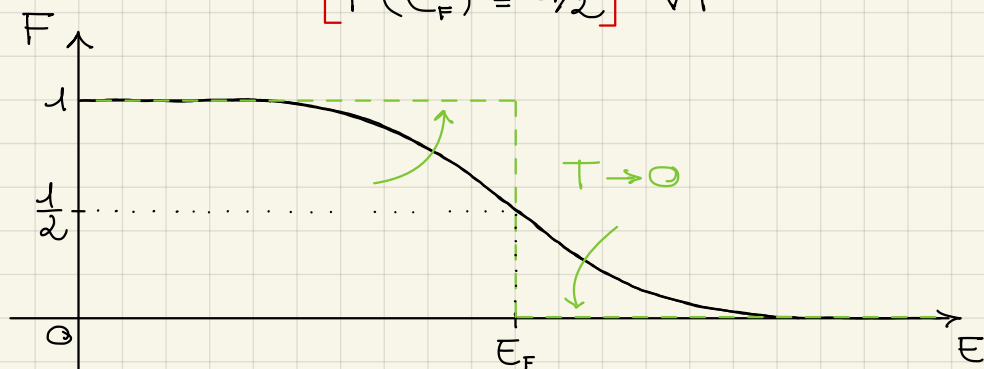
How many of them, however, are actually occupied?

Fermi-Dirac distribution function

$$F(E, T) = \frac{1}{e^{\frac{E-E_F}{kT}} + 1}$$

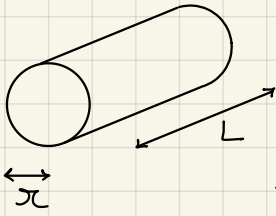
where  $E_F$  is the Fermi level of the material, such that

$$[F(E_F) = 1/2] \quad \forall T$$



$F(E, T)$  is the probability of a state with energy  $E$  at temperature  $T$  of being occupied.

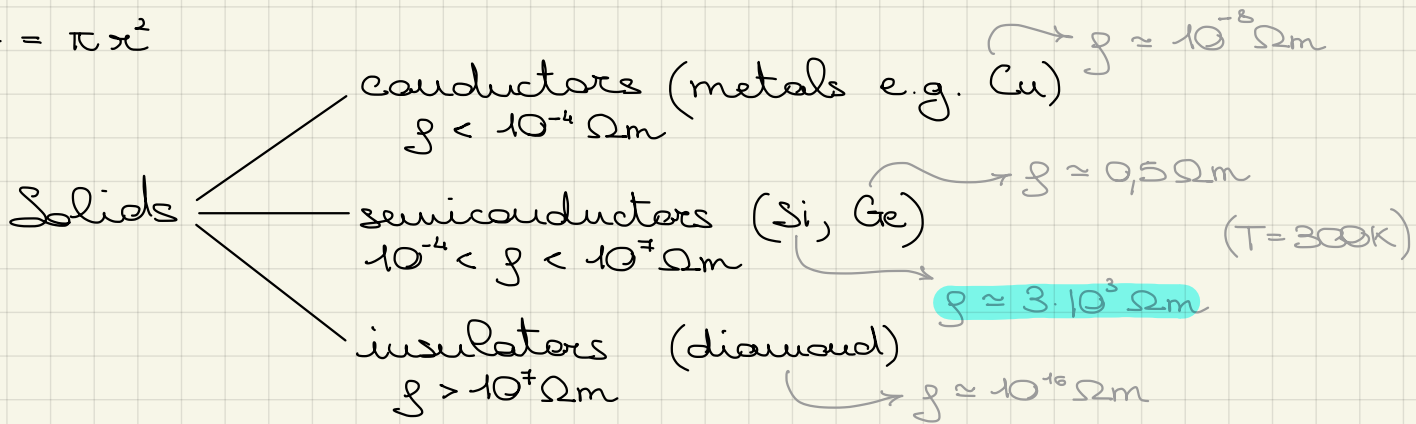
## Conduction in solids



$$\text{resistance } R = \rho \frac{L}{A} = \frac{1}{\sigma} \frac{L}{A}$$

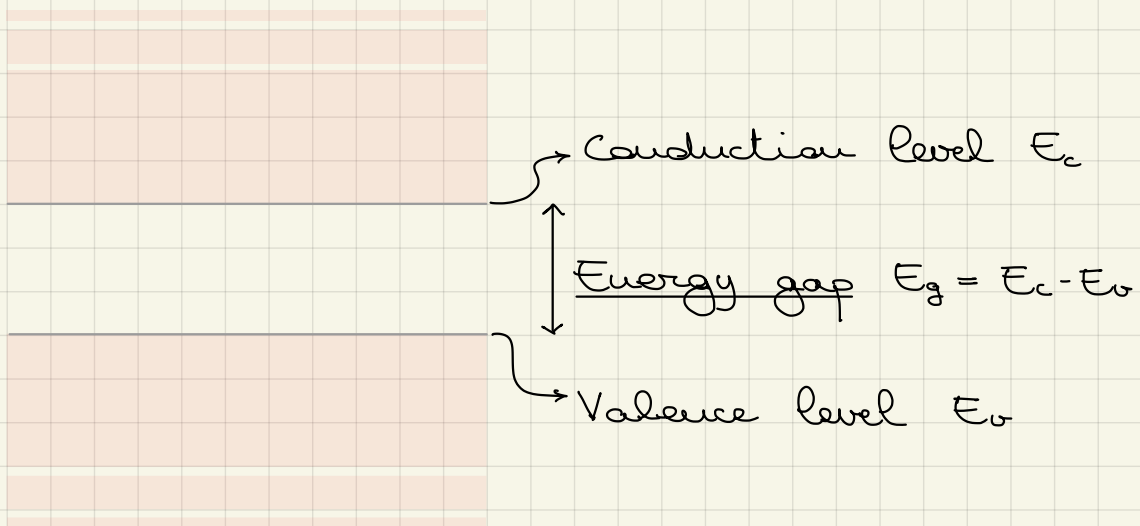
resistivity  $\rho$  [ $\Omega \cdot \text{m}$ ] and conductivity  $\sigma$  [ $\frac{\text{S}}{\text{m}}$ ]

$$A = \pi r^2$$



What determines the conductivity of a material?

From Kronig-Penney theory for energy bands in a crystalline material:



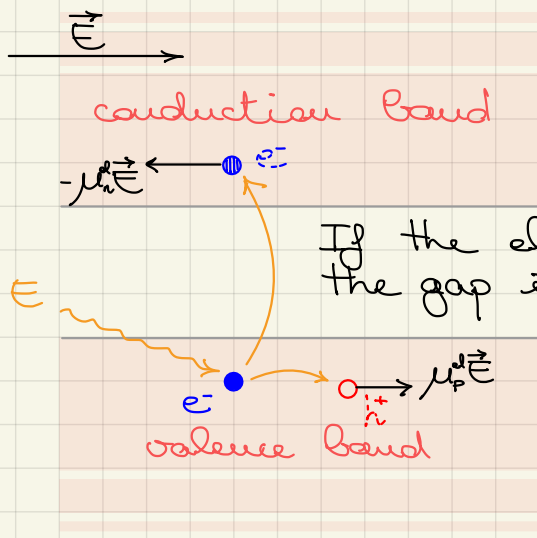
It is the energy gap which determines the conductivity of the materials:

- in conductors (metals), energy distribution is continuous hence  $E_g = 0$
- in semiconductors and insulators,  $E_g$  has a finite non-zero value; the higher  $E_g$ , the less conductive the material

$$E_g^{\text{Si}} = 1,12 \text{ eV}$$

$$E_g^{\text{Ge}} = 0,66$$

$$E_g^{\text{diam}} = 5,47$$



Electrons in the conduction band are free to move in the material, hence they can contribute to conduction

If the electron is provided energy greater than the gap it can undergo "band-to-band" transition

Electrons in the valence band are bound to the covalent bonds of the atoms in the lattice, meaning that they cannot be part of the conduction.

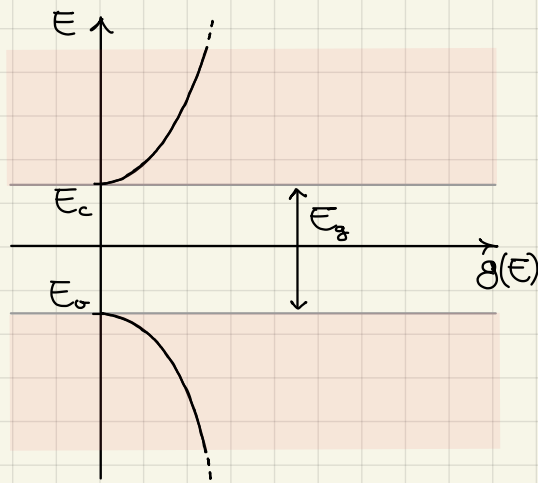
Note that not only electrons but also holes are part of the excitation process. A hole is equivalent to a "negated" electron: its energy is higher the lower it is in the band diagram, hence it naturally occupies higher energy states (where electrons are naturally missing). This also means that holes are conductive in the valence band and they are not in higher energy bands. Exciting an electron therefore produces both a conductive electron in the conductive band and a conductive hole in the valence band.

By applying an external electric field, free electrons and holes will move thus producing a current. The overall current (sum of all contributing electrons and holes) will depend on both the speed of the carriers, which in turn depends on their mobility and on the applied electric field, as well as the number of carriers, which is a function of the energy gap and of the energy provided in the excitation.

Semiconductors are especially interesting since their energy gap is not too high. This means that some electrons in the valence band are able to reach conduction band just by thermal agitation, even at low temperatures.

Between silicon and germanium, the former is preferred thanks to its thermomechanical properties which allow it to endure very high temperatures ( $\sim 1000\text{K}$  for removal of impurities) without melting.

The number of carriers available for conduction can be derived in analogy with the states density and states occupation description for metals:



effective electron/hole mass  
 $1,08 m_0$        $0,81 m_0$

$$C_c = \frac{4\pi (m_n^*)^{3/2}}{h^3}$$

$$C_v = \frac{4\pi (m_p^*)^{3/2}}{h^3}$$

$$g(E) = \begin{cases} g_c(E) = C_c \sqrt{E - E_c} & E \geq E_c \\ g_v(E) = C_v \sqrt{E_v - E} & E \leq E_v \end{cases} \quad \text{where}$$

$$\Rightarrow n = \int_{E_c}^{+\infty} g_c(E) F(E) dE = \frac{N_c}{h^3} e^{-\frac{E_c - E_F}{k_B T}}$$

(# of states)  $\times$  (prob. of occupation)

where

$$N_c = \frac{2}{h^3} (2\pi m_n^* k_B T)^{3/2}$$

$$\Rightarrow p = \int_{-\infty}^{E_v} g_v(E) [1 - F(E)] dE = \frac{N_v}{h^3} e^{-\frac{E_F - E_v}{k_B T}}$$

$$N_v = \frac{2}{h^3} (2\pi m_p^* k_B T)^{3/2}$$

number of conducting electron/holes per unit volume

effective conduction/valence band density of states

From our previous discussion, in pristine conditions (i.e. no doping) it must be  $n = p = n_i$

$$\begin{aligned} N_c e^{-\frac{E_c - E_F}{k_B T}} &= N_v e^{-\frac{E_F - E_v}{k_B T}} \\ e^{\frac{-E_c + E_F + E_F - E_v}{k_B T}} &= \frac{N_c}{N_v} \\ \frac{2E_F - E_c - E_v}{k_B T} &= \ln\left(\frac{N_c}{N_v}\right) \end{aligned}$$

$$E_F = \frac{E_v + E_c}{2} + \frac{k_B T}{2} \ln\left(\frac{N_c}{N_v}\right)$$

neglecting the second term ( $N_c \approx N_v$ )

$$E_F \approx \frac{E_v + E_c}{2}$$

(the Fermi level is about halfway through the gap)

Note that  $n \cdot p = n_i^2 = N_c N_v e^{-\frac{E_c - E_v}{k_B T}} = N_c N_v e^{-\frac{E_g}{k_B T}}$

$$\rightarrow n_i = \sqrt{N_c N_v} e^{-\frac{E_g}{2k_B T}}$$

(T = 300K)

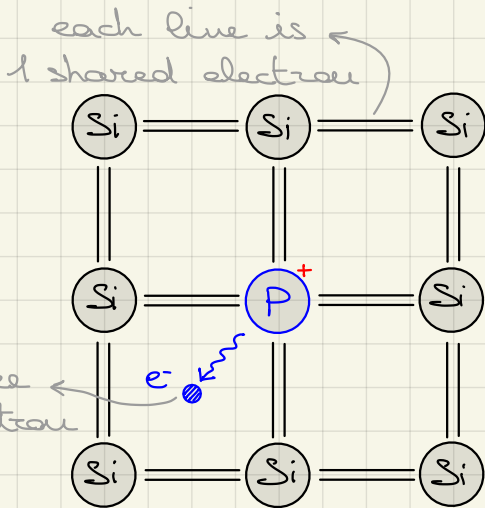
In Si:  $N_c \approx N_v \approx 10^{19} \text{ cm}^{-3}$   $E_g = 1,12 \text{ eV} \rightarrow n_i \approx 10^{10} \text{ cm}^{-3}$

In diamond:  $N_c \approx N_v \approx 10^{19} \text{ cm}^{-3}$   $E_g = 5,47 \text{ eV} \rightarrow n_i \approx 10^{-27} \text{ cm}^{-3}$

It is evident how semiconductors have some (but not many) free carriers at room temperature, unlike insulators which have basically no conducting carrier whatsoever.

It is possible to exploit the below-average conductivity and resilience to thermomechanical stress of semiconductors (especially silicon) to realize new materials whose resistance can be accurately set and modulated.

### Doping



The process of doping consists of the substitution of one silicon atom in the lattice with an atom from the 3rd or 5th group of the Periodic Table, typically boron or phosphorus. Note that silicon belongs to the 4th group.

The result is that the doping atom will have either too few or too many electrons to perform the four bonds of the tetrahedron. As a consequence, the dopant will have to either steal or release an electron thus increasing the number of free holes or free electrons. Since charge neutrality has to be maintained, the dopant will negatively or positively ionize.

Dopants of the former type are called acceptors, while of the latter are called donors.

Since we can decide how many doping atoms we introduce in the lattice, we can also control the concentration of conducting carriers:

$$p = p(N_A^-, N_0^+) \quad n = n(N_0^+, N_A^-)$$

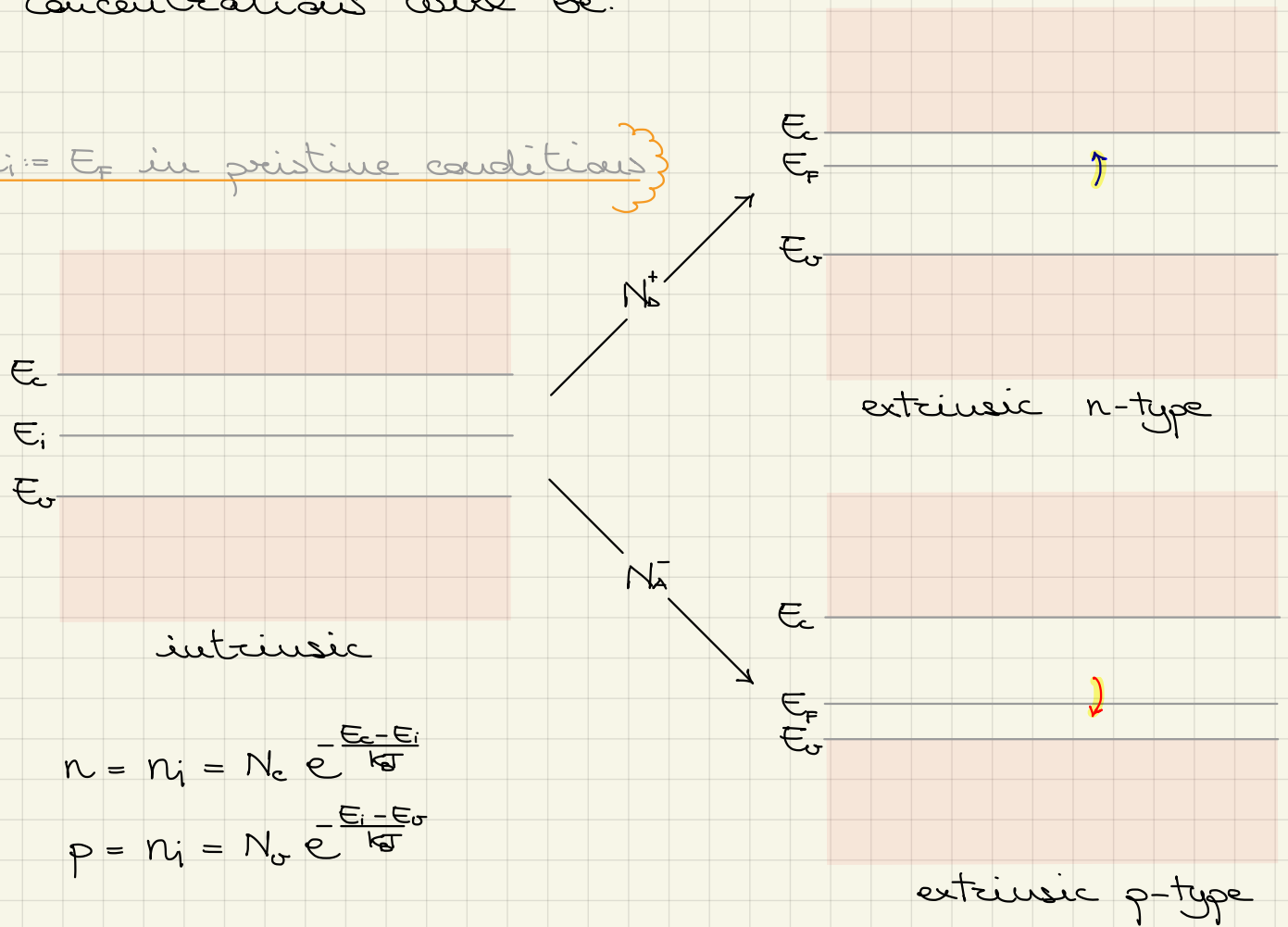
# of acceptors ← → # of donors



We are modifying the concentration from intrinsic to extrinsic.

Let's see what this means from the energetic standpoint, and what the resulting carrier concentrations will be:

$E_i = E_F$  in pristine conditions



$$n = n_i = N_c e^{-\frac{E_c - E_i}{k_B T}}$$

$$p = n_i = N_v e^{-\frac{E_i - E_v}{k_B T}}$$

"Doping moves the Fermi energy level"

$$n = N_c e^{-\frac{E_c - E_F}{k_B T}}$$

$$p = N_v e^{-\frac{E_F - E_v}{k_B T}}$$

$$\Rightarrow \left[ n = n_i e^{\frac{E_F - E_i}{k_B T}} \quad p = n_i e^{\frac{E_i - E_F}{k_B T}} \right]$$

$$\Rightarrow \boxed{p \cdot n = n_i^2} \quad \text{Mass-action law} \quad (\text{always valid at equilibrium})$$

Another consequence is:

$$\left[ \text{Thermal equilibrium} \iff E_F \text{ is a constant} \right]$$

promotion and demotion of carriers happen at the same rate



$$\begin{cases} \text{Mass-action law} \rightarrow n \cdot p = n_i^2 \rightarrow np - n_i^2 = 0 \\ \text{Electroneutrality} \rightarrow \rho^e = 0 \rightarrow q(\underbrace{N_D^+ - N_A^-}_{\text{fixed charges}} + \underbrace{p - n}_{\text{mobile charges}}) = 0 \end{cases}$$

Let us define  $\phi := N_D^+ - N_A^-$

We will now suppose to be dealing with a n-type region (i.e.  $\phi \gg 0$ ):

$$\begin{cases} p = \frac{n_i^2}{n} \\ \phi + p - n = 0 \end{cases} \Rightarrow n\phi + n_i^2 - n^2 = 0 \rightarrow \left[ n = \frac{\phi + \sqrt{\phi^2 + 4n_i^2}}{2} \right]$$

$N_D^+ \gg n_i$

$$\left[ E_F = E_i + k_B T \ln \frac{N_D^+}{n_i} \right] \leftarrow \begin{matrix} N_D^+ = n_i e^{\frac{E_F - E_i}{k_B T}} \\ \frac{n_i^2}{N_D^+} = n_i e^{\frac{E_i - E_F}{k_B T}} \end{matrix} \leftarrow \begin{matrix} n \approx N_D^+ \\ p \approx \frac{n_i^2}{N_D^+} \end{matrix}$$

By analogy, in a p-type region ( $\phi \ll 0$ ):

$$\left[ p = \frac{-\phi + \sqrt{\phi^2 + 4n_i^2}}{2} \right]$$

$N_A^- \gg n_i$

$$\begin{matrix} p \approx N_A^- \\ n \approx \frac{n_i^2}{N_A^-} \end{matrix} \Rightarrow \begin{matrix} \frac{n_i^2}{N_A^-} = n_i e^{\frac{E_F - E_i}{k_B T}} \\ N_A^- = n_i e^{\frac{E_i - E_F}{k_B T}} \end{matrix} \Rightarrow \left[ E_F = E_i - k_B T \ln \frac{N_A^-}{n_i} \right]$$

In non-equilibrium conditions, this formulas for the Fermi level do not hold anymore and a more complex model from thermodynamics should be adopted.

In practice, effective values for the Fermi level, called "quasi-Fermi levels"  $E_{Fn}$  and  $E_{Fp}$ , are used within all the above theory, instead of adopting a whole new model.

$$\left[ n = n_i e^{\frac{E_{Fn} - E_i}{k_B T}} \quad p = n_i e^{\frac{E_i - E_{Fp}}{k_B T}} \right]$$

$$\left[ n \cdot p = n_i^2 e^{\frac{E_{Fn} - E_{Fp}}{k_B T}} \right]$$

In non-equilibrium:  $E_{Fn} \neq E_{Fp}$ . If  $E_{Fn} > E_{Fp}$  then  $np \gg n_i^2$  which means that carrier concentration is overwhelming and recombination predominant, thus restoring equilibrium. Viceversa, if  $E_{Fn} < E_{Fp}$  then carrier concentration is

underwhelming and generation predominant, thus again restoring equilibrium.

## Drift-diffusion model

3 balance equations

$$\begin{cases} \vec{\nabla} \cdot \vec{D} = g^{\text{ext}} \\ -q \frac{\partial n}{\partial t} + \vec{\nabla} \cdot \vec{J}_n = q(R-G) \\ q \frac{\partial p}{\partial t} + \vec{\nabla} \cdot \vec{J}_p = -q(R-G) \end{cases} \quad \left. \begin{array}{l} \text{Ampère - Maxwell} \\ \text{law} \\ \vec{\nabla} \cdot (\vec{J} + \frac{\partial \vec{D}}{\partial t}) = 0 \end{array} \right\}$$

$$g^{\text{ext}} = q(N_D^+ - N_A^-) + q(p-n) \quad \vec{D} = \epsilon \vec{E} = \epsilon_0 \epsilon_r \vec{E}$$

$$\vec{J}_n = q \mu_n^{\text{eff}} n \vec{E} + q D_n \vec{\nabla} n$$

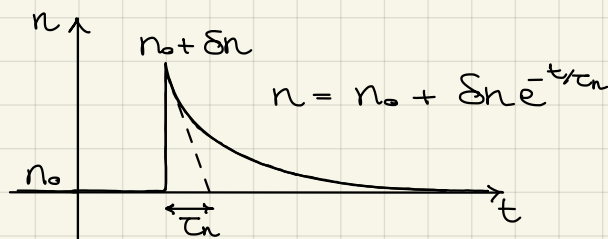
$$\vec{J}_p = q \mu_p^{\text{eff}} p \vec{E} - q D_p \vec{\nabla} p$$

Shockley-Read-Hall recombination formula

$$R - G = \frac{p n - n_i^2}{\tau_n(p + n_i) + \tau_p(n + n_i)} = 0 \iff \text{Therm. equ. (np = n_i^2)}$$

carrier lifetimes

measure the rapidity with which an excited carrier is reabsorbed back to neutrality

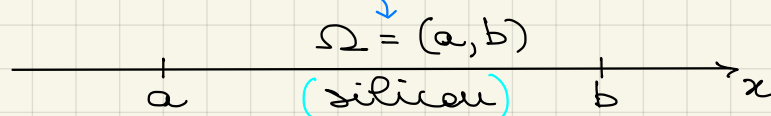


$\tau_n$  and  $\tau_p$  strongly depend on the atomic structure of the material and on the presence of dopants (hence they might be non-uniform in the domain).

$$\vec{E} = -\vec{\nabla} \psi$$

voltage

Simplifications: 1D + stationary  $\rightarrow \frac{\partial n}{\partial t} = \frac{\partial p}{\partial t} = 0$



Putting everything together:

$$(1) \begin{cases} \frac{\partial}{\partial x} (-\varepsilon \frac{\partial \Psi}{\partial x}) - q_p + q_n - q \phi = 0 \\ -\frac{\partial}{\partial x} (-q \mu_n^{\text{el}} n \frac{\partial \Psi}{\partial x} + q D_n \frac{\partial n}{\partial x}) + q \frac{p \cdot n - n_i^2}{\tau_n(p+n_i) + \tau_p(n+n_i)} = 0 \\ \frac{\partial}{\partial x} (-q \mu_p^{\text{el}} p \frac{\partial \Psi}{\partial x} - q D_p \frac{\partial p}{\partial x}) + q \frac{p \cdot n - n_i^2}{\tau_n(p+n_i) + \tau_p(n+n_i)} = 0 \end{cases} \quad x \in \Omega$$

$$F(\underline{u}) = \underline{0}$$

we assume it, for now, to be unique

solve as a fixed point iteration (like in a cable equation problem):

find  $\underline{u}^*$  such that

$$F(\underline{u}^*) = \underline{0}$$

becomes

find  $\underline{u}^*$  such that

$$\underline{u}^* = T_F(\underline{u}^*)$$

Let's first explicitly write what are  $\underline{u}$  and  $F$  of problem (1):

$$\underline{u} = \begin{bmatrix} \psi(x) \\ n(x) \\ p(x) \end{bmatrix} \quad F(\underline{u}) = \begin{bmatrix} f_\psi(\underline{u}) \\ f_n(\underline{u}) \\ f_p(\underline{u}) \end{bmatrix}$$

- $f_\psi(\underline{u}) = \frac{\partial}{\partial x} (-\varepsilon \frac{\partial \Psi}{\partial x}) - q_p + q_n - q \phi$
- $f_n(\underline{u}) = -\frac{\partial}{\partial x} (-q \mu_n^{\text{el}} n \frac{\partial \Psi}{\partial x} + q D_n \frac{\partial n}{\partial x}) + q \mathcal{R}(n, p)$
- $f_p(\underline{u}) = \frac{\partial}{\partial x} (-q \mu_p^{\text{el}} p \frac{\partial \Psi}{\partial x} - q D_p \frac{\partial p}{\partial x}) + q \mathcal{R}(n, p)$

$$\text{where } \mathcal{R}(n, p) = \frac{p \cdot n - n_i^2}{\tau_n(p+n_i) + \tau_p(n+n_i)}$$

Now, what fixed point operator  $T_F$  should we use?

1) Newton's method

Given  $\underline{u}^{(0)} = \begin{bmatrix} \psi^{(0)}(x) \\ n^{(0)}(x) \\ p^{(0)}(x) \end{bmatrix}$ ,  $\forall k \geq 0$  until convergence:

Fréchet derivative

$$\text{solve: } \underline{F}'(\underline{u}^{(k)}) \cdot \underline{\delta u}^{(k)} = -F(\underline{u}^{(k)})$$

residual

where  $\underline{F}'$  is the Jacobian of  $F$  such that  $(\underline{F}')_{i,j} = \frac{\partial F_i}{\partial u_j}$

$$\text{update: } \underline{u}^{(k+1)} = \underline{u}^{(k)} + \underline{\delta u}^{(k)}$$

increment

As already seen, Newton's method benefits from local convergence of second order. In other words:

$$\exists B^* \ni \underline{u}^* \text{ such that } \forall \underline{u}^{(0)} \in B^*: \lim_{k \rightarrow \infty} \underline{u}^{(k)} = \underline{u}^*$$

$$\text{and also } \|\underline{u}^{(k+1)} - \underline{u}^*\|_V \leq C \|\underline{u}^{(k)} - \underline{u}^*\|_V^2 \quad \forall k \geq k_0 > 0$$

i.e. for  $k$  large enough  $B^* \subseteq V$



This is a nice property, which however requires the initial guess  $\underline{u}^{(0)}$  to be sufficiently close to  $\underline{u}^*$  (i.e. within  $B^*$ ) to guarantee not only convergence order, but also convergence itself.

→ Newton's method has very tight requirements to work properly, but when they are granted it is one of the fastest and most reliable among all fixed point methods.

Let's now determine the explicit form of all terms needed to run Newton's iteration:

$$(2) \quad \begin{array}{|c|c|c|c|} \hline \frac{\partial f_\psi}{\partial \phi}(\underline{u}^{(k)}) & \frac{\partial f_\psi}{\partial n}(\underline{u}^{(k)}) & \frac{\partial f_\psi}{\partial p}(\underline{u}^{(k)}) & \delta \psi^{(k)} \\ \hline \frac{\partial f_n}{\partial \phi}(\underline{u}^{(k)}) & \frac{\partial f_n}{\partial n}(\underline{u}^{(k)}) & \frac{\partial f_n}{\partial p}(\underline{u}^{(k)}) & \delta n^{(k)} \\ \hline \frac{\partial f_p}{\partial \phi}(\underline{u}^{(k)}) & \frac{\partial f_p}{\partial n}(\underline{u}^{(k)}) & \frac{\partial f_p}{\partial p}(\underline{u}^{(k)}) & \delta p^{(k)} \\ \hline \end{array} = - \begin{array}{|c|} \hline f_\psi(\underline{u}^{(k)}) \\ \hline f_n(\underline{u}^{(k)}) \\ \hline f_p(\underline{u}^{(k)}) \\ \hline \end{array}$$

$$\underline{F}'(\underline{u}^{(k)}) \cdot \underline{\delta u}^{(k)} = -F(\underline{u}^{(k)})$$

$$\bullet \frac{\partial f_\psi}{\partial \phi}(\underline{u}^{(k)}) = \frac{\partial}{\partial x} \left( -\varepsilon \frac{\partial}{\partial x} (\cdot) \right)$$

$$\bullet \frac{\partial f_\psi}{\partial n}(\underline{u}^{(k)}) = q(\cdot)$$

- $\frac{\partial \mathcal{L}}{\partial \psi}(\underline{u}^{(k)}) = -q(\cdot)$
- $\frac{\partial \mathcal{L}}{\partial \phi}(\underline{u}^{(k)}) = -\frac{\partial}{\partial x}(-q \mu_n^{el} n^{(k)} \frac{\partial}{\partial x}(\cdot))$
- $\frac{\partial \mathcal{L}}{\partial n}(\underline{u}^{(k)}) = -\frac{\partial}{\partial x}(-q \mu_n^{el}(\cdot) \frac{\partial \psi^{(k)}}{\partial x} + q D_n \frac{\partial}{\partial x}(\cdot)) + q \frac{\partial \mathcal{R}}{\partial n}(n^{(k)}, p^{(k)})(\cdot)$
- $\frac{\partial \mathcal{L}}{\partial p}(\underline{u}^{(k)}) = q \frac{\partial \mathcal{R}}{\partial p}(n^{(k)}, p^{(k)})(\cdot)$
- $\frac{\partial \mathcal{L}}{\partial \phi}(\underline{u}^{(k)}) = \frac{\partial}{\partial x}(-q \mu_p^{el} p^{(k)} \frac{\partial}{\partial x}(\cdot))$
- $\frac{\partial \mathcal{L}}{\partial n}(\underline{u}^{(k)}) = q \frac{\partial \mathcal{R}}{\partial n}(n^{(k)}, p^{(k)})(\cdot)$
- $\frac{\partial \mathcal{L}}{\partial p}(\underline{u}^{(k)}) = \frac{\partial}{\partial x}(-q \mu_p^{el}(\cdot) \frac{\partial \psi^{(k)}}{\partial x} - q D_p \frac{\partial}{\partial x}(\cdot)) + q \frac{\partial \mathcal{R}}{\partial p}(n^{(k)}, p^{(k)})(\cdot)$

Notes:

- on the diagonal, all terms are associated to the variation (derivative in  $x$ ) of their respective quantity ( $\delta u_i^{(k)}$ )
- in the first column (associated to  $\delta \psi^{(k)}$ ), all terms have the form of an elliptic equation

We now have to make the problem computable from a numerical standpoint (finite element method for differential equations of each iteration).



Boundary conditions:

	$\psi(0) = \bar{\psi}_0$	$\psi(L) = \bar{\psi}_L$
↓ Dirichlet	$n(0) = \bar{n}_0 > 0$	$n(L) = \bar{n}_L > 0$
	$p(0) = \bar{p}_0 > 0$	$p(L) = \bar{p}_L > 0$

At the boundaries it has to be:  $\delta \psi(0) = \delta \psi(L) = 0$   
homogeneous Dirichlet B.c.  $\left\{ \begin{array}{l} \delta n(0) = \delta n(L) = 0 \\ \delta p(0) = \delta p(L) = 0 \end{array} \right.$

since  $\underline{u}^{(k)}(0) = \underline{u}_0$  and  $\underline{u}^{(k)}(L) = \underline{u}_L \quad \forall k$ , therefore:

$$\delta\psi \in H_0^1(\Omega)$$

$$\delta n \in H_0^1(\Omega)$$

$$\delta p \in H_0^1(\Omega)$$

Weak formulation of problem (2):

$$\left\{ \begin{aligned} \int_0^L \left[ \frac{\partial}{\partial x} (-\varepsilon \frac{\partial}{\partial x} (\delta\psi)) + q \delta n - q \delta p \right] \phi^i &= - \int_0^L j_\psi \phi^i \quad \forall \phi^i \in H_0^1(\Omega) \\ \int_0^L \left[ -\frac{\partial}{\partial x} (-q \mu_n^d n \frac{\partial}{\partial x} (\delta\psi)) - \frac{\partial}{\partial x} (-q \mu_n^d \delta n \frac{\partial \psi}{\partial x} + q D_n \frac{\partial}{\partial x} (\delta n)) + q \frac{\partial R}{\partial n} \delta n + q \frac{\partial R}{\partial p} \delta p \right] \phi^{ii} &= - \int_0^L j_n \phi^{ii} \quad \forall \phi^{ii} \in H_0^1(\Omega) \\ \int_0^L \left[ \frac{\partial}{\partial x} (-q \mu_p^d p \frac{\partial}{\partial x} (\delta\psi)) + q \frac{\partial R}{\partial n} \delta n + \frac{\partial}{\partial x} (-q \mu_p^d \delta p \frac{\partial \psi}{\partial x} - q D_p \frac{\partial}{\partial x} (\delta p)) + q \frac{\partial R}{\partial p} \delta p \right] \phi^{iii} &= - \int_0^L j_p \phi^{iii} \quad \forall \phi^{iii} \in H_0^1(\Omega) \end{aligned} \right.$$

$$\int_0^L \varepsilon \frac{\partial}{\partial x} (\delta\psi) \cdot \frac{\partial \phi^i}{\partial x} + \int_0^L q \delta n \phi^i - \int_0^L q \delta p \phi^i = - \int_0^L j_\psi \phi^i$$

$$- \int_0^L q \mu_n^d n \frac{\partial}{\partial x} (\delta\psi) \frac{\partial \phi^{ii}}{\partial x} + \int_0^L (-q \mu_n^d \delta n \frac{\partial \psi}{\partial x} + q D_n \frac{\partial}{\partial x} (\delta n)) \frac{\partial \phi^{ii}}{\partial x} + \int_0^L q \frac{\partial R}{\partial n} \delta n \phi^{ii} + \int_0^L q \frac{\partial R}{\partial p} \delta p \phi^{ii} = - \int_0^L j_n \phi^{ii}$$

$$\int_0^L q \mu_p^d p \frac{\partial}{\partial x} (\delta\psi) \frac{\partial \phi^{iii}}{\partial x} + \int_0^L q \frac{\partial R}{\partial n} \delta n \phi^{iii} + \int_0^L (q \mu_p^d \delta p \frac{\partial \psi}{\partial x} + q D_p \frac{\partial}{\partial x} (\delta p)) \frac{\partial \phi^{iii}}{\partial x} + \int_0^L q \frac{\partial R}{\partial p} \delta p \phi^{iii} = - \int_0^L j_p \phi^{iii}$$

Space discretization:  $\delta u \rightsquigarrow \delta u_n$



$$\delta\psi_n(x) = \sum_{j=1}^{M_n+1} \delta\psi_j \cdot \phi_j(x) = \sum_{j=2}^{M_n} \delta\psi_j \cdot \phi_j(x)$$

↑  
homogeneous b.c.

← basis functions

Repeating the same process for  $\delta n$  and  $\delta p$ :

$$\left\{ \begin{aligned} \underline{K}_{\psi\psi} \delta\psi + \underline{K}_{\psi n} \delta n + \underline{K}_{\psi p} \delta p &= -\underline{F}_\psi \\ \underline{K}_{n\psi} \delta\psi + \underline{K}_{nn} \delta n + \underline{K}_{np} \delta p &= -\underline{F}_n \\ \underline{K}_{p\psi} \delta\psi + \underline{K}_{pn} \delta n + \underline{K}_{pp} \delta p &= -\underline{F}_p \end{aligned} \right.$$



- Notes:
- $\underline{K}_{u_i, u_i}$  has all elements of the first and last row equal to 0, except the first and last element, respectively, which are equal to 1
  - $\underline{K}_{u_i, u_j} \forall i \neq j$  has all elements of the first and last row equal to 0
  - $\underline{F}_{u_i}$  has the first and last element equal to 0
- These are necessary conditions for enforcing homogeneous B.C.

A few additional comments:

- We considered, throughout the entirety of the discussion,  $\mu^{\text{eff}}$  to be constant; however, for large electric fields (big  $\psi$ ), carrier velocity saturates meaning that the mobility is not constant but depends on  $\psi$  itself. This physical description of  $\mu^{\text{eff}}$  would make our model more accurate but also much more complicated.
- Matrices we obtained for the finite element method can be badly conditioned because of the several diverse units and orders of magnitude involved. A technique called "balancing" is typically adopted in numerical solvers to improve the matrix condition number.
- Convergence speed of the method can be hindered by the non-exact description of the derivatives that are contained in the iterative matrices

## 2) Gummel's method

Given  $\underline{x}^{(0)}$ ,  $\forall k \geq 0$  until convergence:

$$\underline{x}^{(k+1)} = \underline{T}_G(\underline{x}^{(k)})$$

For Newton's method it was:  $\underline{u}^{(k+1)} = \underline{T}_N(\underline{u}^{(k)})$  ( $\underline{u} \neq \underline{x}$  !!!)

$$\underline{T}_N(\underline{u}) = \underline{u} - (\underline{F}'(\underline{u}))^{-1} \underline{F}(\underline{u})$$



Let's see what  $\underline{x}$  and  $T_0$  stand for.

$$\underline{x} = \begin{bmatrix} \phi_n \\ \phi_p \end{bmatrix}$$

where  $\phi_n$  and  $\phi_p$  embody the quasi-Fermi potentials of  $n$  and  $p$ , respectively.

$\textcircled{a} \quad n = n_i e^{\frac{\psi - \phi_n}{V_{th}}}$       $\textcircled{b} \quad p = n_i e^{\frac{\phi_p - \psi}{V_{th}}}$ 
→ from Maxwell-Boltzmann statistics

Let us now introduce the **non-linear** operators:

- $\textcircled{1} \quad \mathcal{N}_\psi(\psi, \underline{x}) = 0 \rightarrow$  non-linear w.r.t.  $\psi$
- $\textcircled{2} \quad \mathcal{N}_n(n, \psi(\underline{x}), p(\underline{x}, \psi(\underline{x}))) = 0 \rightarrow$  non-linear w.r.t.  $n$
- $\textcircled{3} \quad \mathcal{N}_p(p, \psi(\underline{x}), n(\underline{x}, \psi(\underline{x}))) = 0 \rightarrow$  non-linear w.r.t.  $p$

which written in full are:

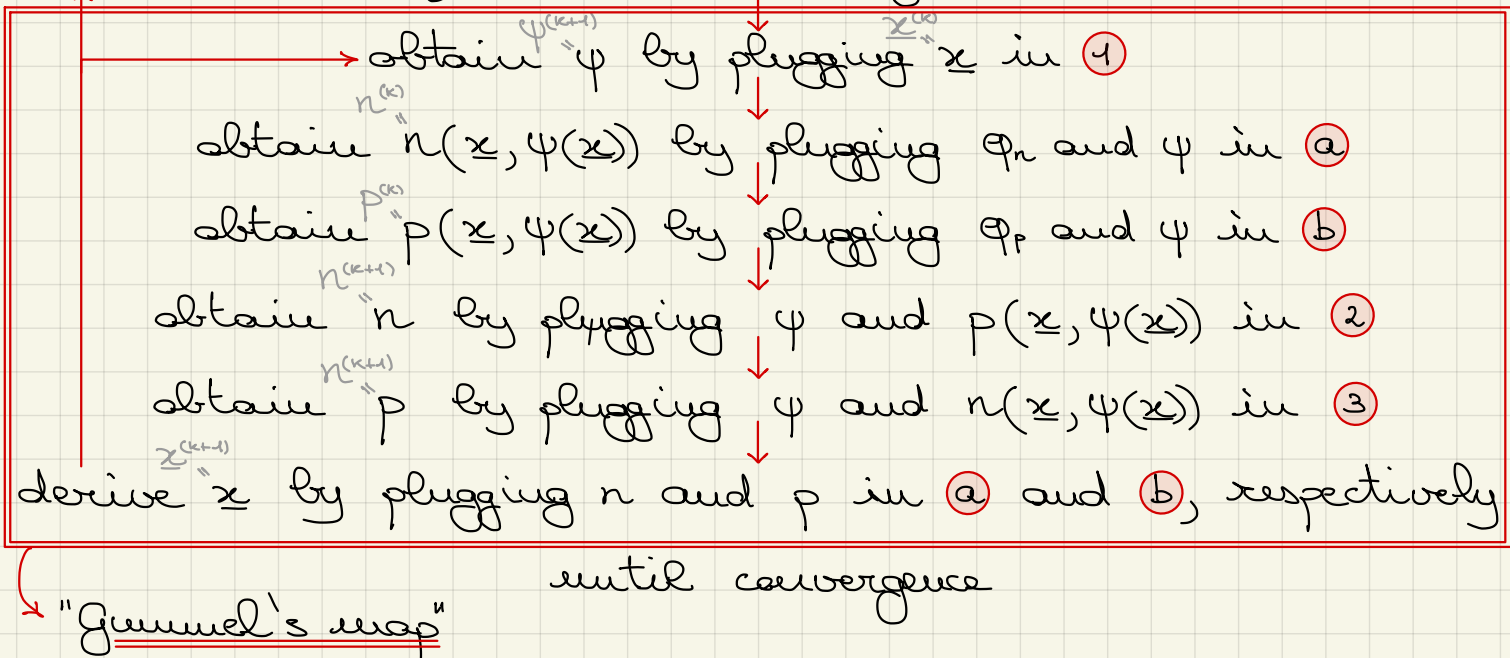
$$\mathcal{N}_\psi = \frac{\partial}{\partial x} \left( -\epsilon \frac{\partial \psi}{\partial x} \right) + q n_i e^{\frac{\psi - \phi_n}{V_{th}}} - q n_i e^{\frac{\phi_p - \psi}{V_{th}}} - q \mathcal{D}$$

$$\mathcal{N}_n = -\frac{\partial}{\partial x} \left[ -q \mu_n^{\text{eff}} n \frac{\partial \psi(\underline{x})}{\partial x} + q D_n \frac{\partial n}{\partial x} \right] + q \frac{p(\underline{x}, \psi(\underline{x})) \cdot n - n_i^2}{\tau_n(p(\underline{x}, \psi(\underline{x})) + n_i) + \tau_p(n(\underline{x}, \psi(\underline{x})) + n_i)}$$

$$\mathcal{N}_p = \frac{\partial}{\partial x} \left[ -q \mu_p^{\text{eff}} p \frac{\partial \psi(\underline{x})}{\partial x} - q D_p \frac{\partial p}{\partial x} \right] + q \frac{p \cdot n(\underline{x}, \psi(\underline{x})) - n_i^2}{\tau_n(p + n_i) + \tau_p(n(\underline{x}, \psi(\underline{x})) + n_i)}$$

Then, the approach of Gummel's method is to:

$\underline{x}^{(k+1)}$  given the initial guess  $\underline{x}^{(0)}$



Therefore:  $T_G(\underline{x}^{(k)}) = \begin{bmatrix} \psi^{(k+1)} - \sqrt{H} \ln\left(\frac{n^{(k+1)}}{n_i}\right) \\ \psi^{(k+1)} + \sqrt{H} \ln\left(\frac{p^{(k+1)}}{n_i}\right) \end{bmatrix}$

$\rightarrow$  inverse of (a)  
 $\rightarrow$  inverse of (b)

where  $\psi^{(k+1)}$ ,  $n^{(k+1)}$  and  $p^{(k+1)}$  are the ones obtained in Gummel's map through the non-linear Poisson and continuity equations.

Issue: computation of non-linear equations.

A possible solution to partially overcome the non-linear nature of the method, which otherwise would be extremely simple and straightforward, is to use the "lagging" technique for what concerns the computation of  $n^{(k+1)}$  and  $p^{(k+1)}$  (i.e. for equations (2) and (3))

"lagging" means to substitute the unknown values to be computed for the current step ( $n^{(k+1)}$  and  $p^{(k+1)}$ ) that make up the non-linearity with the known values already computed for the previous step ( $n^{(k)}$  and  $p^{(k)}$ ). By applying this concept to our equations we get:

$$\tilde{N}_n = -\frac{\partial}{\partial x} \left[ -q \mu_n^{eff} n \frac{\partial \psi(x)}{\partial x} + q D_n \frac{\partial n}{\partial x} \right] + q \frac{p(x, \psi(x)) n - n_i^2}{\tau_n(p(x, \psi(x)) + n_i) + \tau_p(n(x, \psi(x)) + n_i)}$$

*linear terms are untouched*

$$\tilde{N}_p = \frac{\partial}{\partial x} \left[ -q \mu_p^{eff} p \frac{\partial \psi(x)}{\partial x} - q D_p \frac{\partial p}{\partial x} \right] + q \frac{p \cdot n(x, \psi(x)) - n_i^2}{\tau_n(p(x, \psi(x)) + n_i) + \tau_p(n(x, \psi(x)) + n_i)}$$

*non-linear terms are removed*

which are linear advection-diffusion-reaction equations thus inheriting all properties associated with them (such as the maximum principle, as we will see).

Note that this approach would not work with (1) since we don't have a value of  $\psi^{(k)}$  to substitute in the non-linearity.

$$\begin{cases} \frac{\partial}{\partial x} \left( -\varepsilon \frac{\partial \psi}{\partial x} \right) + q n_i e^{\frac{\psi - \phi_n}{\sqrt{H}}} - q n_i e^{\frac{\phi_p - \psi}{\sqrt{H}}} - q \phi = 0, & x \in (0, L) \\ \psi(0) = \bar{\psi}_0, & \psi(L) = \bar{\psi}_L \end{cases}$$

To overcome also this non-linearity, we could use Newton's method:

given  $\psi^{(0)}$ ,  $\forall j \geq 0$  compute  $\{\psi^{(j)}\}$  until convergence

solve:  $\mathcal{N}'_{\psi}(\psi^{(j)}) \delta\psi^{(j)} = -\mathcal{N}_{\psi}(\psi^{(j)})$

update:  $\psi^{(j+1)} = \psi^{(j)} + \delta\psi^{(j)}$

where  $\mathcal{N}'_{\psi}(\psi^{(j)}) \cdot \delta\psi^{(j)} = \frac{\partial}{\partial x} \left( -\varepsilon \frac{\partial}{\partial x} (\delta\psi^{(j)}) \right) + q n_i \left[ \frac{e^{\frac{\varphi_p - \psi^{(j)}}{V_{th}}}}{V_{th}} + \frac{e^{\frac{\psi^{(j)} - \varphi_n}{V_{th}}}}{V_{th}} \right] \delta\psi^{(j)}$

Also this equation, to be solved for each step of Newton's iteration, is an advection-diffusion-reaction problem.

Existence, Uniqueness and convergence of a solution of Gummel's map

Uniqueness: we will just assume that  $\exists!$  solution of the problem denoted as

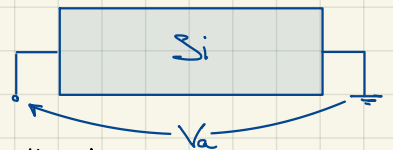
$\psi^*, \varphi_n^*, \varphi_p^*, n^* = n_i e^{\frac{\varphi_p^* - \psi^*}{V_{th}}}, p^* = n_i e^{\frac{\psi^* - \varphi_n^*}{V_{th}}}$

Existence:

introduce  $V_{\psi} := \{ \psi \in L^2(\Omega) \mid \alpha \leq \psi(x) \leq \beta \quad \forall x \in \bar{\Omega} \}$

and suppose  $\psi^{(0)} \in V_{\psi}$ .

$\alpha = \min(V_a, 0)$   
 $\beta = \max(V_a, 0)$



By these suppositions it can be shown that:

$\implies \forall k \geq 0 \quad \psi(\psi^{(k)}) = \psi^{(k)} \in V_{\psi}$

where  $V_{\psi} := \{ w \in H^1(\Omega) \mid A \leq w(x) \leq B \quad \forall x \in \Omega \}$

$A = \min(\check{\psi}_0, \check{\psi})$       $\check{\psi}_0 = \min(\check{\psi}_0, \check{\psi}_L)$   
 $B = \max(\check{\psi}_0, \check{\psi})$       $\check{\psi}_0 = \max(\check{\psi}_0, \check{\psi}_L)$

*depend on the doping*

$\implies \forall k \geq 0 \quad \psi^{(k)} \in V_{\psi}$

$V_{\psi}$  is called "invariant region" for the quasi-Fermi potentials.

$V_{\psi}$  is an "invariant region" for the potential.

This grants the existence of the fixed point of the iteration since

$V_{\psi} \rightarrow T_G(V_{\psi}) \subset V_{\psi}$

Convergence: it can be demonstrated that the contraction condition

$$\exists c < 1 \text{ such that } \|y^{(k)} - x^{(k)}\|_{(L^2(\Omega))^2} \leq c \|y^{(k-1)} - x^{(k-1)}\|_{(L^2(\Omega))^2}$$

$$\forall x^{(k)}, y^{(k)} \in V_x \quad x^{(k)} \neq y^{(k)} \quad x^{(k)} = T_G(x^{(k-1)}) \quad y^{(k)} = T_G(y^{(k-1)})$$

which grants convergence for  $T_G \forall x^{(k)} \in V_x$ , holds true for  $\forall \epsilon$  small enough (i.e.  $\forall \epsilon$  comparable with  $V_{th}$ )

Well-posedness of the current continuity equations

$$\textcircled{2} \quad \tilde{X}_n(n) = 0$$

$$\textcircled{3} \quad \tilde{X}_p(p) = 0$$

as already pointed out, can be written in the form:

$$\frac{\partial J(u)}{\partial x} + \alpha \cdot u = g$$

$$J(u) = V \cdot u - D \frac{\partial u}{\partial x}$$

where in case  $\textcircled{2}$  it is:

$$u = n^{(k+1)} \quad V = -\mu_n^{el} E^{(k+1)} \quad D = D_n$$

$$\alpha = \frac{p^{(k)}}{\tau_n(p^{(k)} + n_i) + \tau_p(n^{(k)} + n_i)} > 0 \quad g = \frac{n_i^2}{\tau_n(p^{(k)} + n_i) + \tau_p(n^{(k)} + n_i)} > 0$$

while in case  $\textcircled{3}$  it is:

$$u = p^{(k+1)} \quad V = \mu_p^{el} E^{(k+1)} \quad D = D_p$$

$$\alpha = \frac{n^{(k)}}{\tau_n(p^{(k)} + n_i) + \tau_p(n^{(k)} + n_i)} > 0 \quad g = \frac{n_i^2}{\tau_n(p^{(k)} + n_i) + \tau_p(n^{(k)} + n_i)} > 0$$

(remembering that  $E^{(k+1)} = -\frac{\partial \psi^{(k+1)}}{\partial x}$  is a given quantity).

Dirichlet B.c.:  $u(0) = \bar{u}_0 > 0$  and  $u(L) = \bar{u}_L > 0$  are also to be enforced to solve the two equations.

To demonstrate the validity of maximum principle, we apply a notation change to the unknowns  $n$  and  $p$  that comes from equations  $\textcircled{a}$  and  $\textcircled{b}$ .

$$n = g_n e^{\frac{\psi}{V_{th}}}$$

where

$$p = g_p e^{-\frac{\psi}{V_{th}}}$$

$$g_n = n_i e^{-\frac{\phi}{V_{th}}}$$

$$g_p = n_i e^{\frac{\phi}{V_{th}}}$$

effective concentrations

Slotboom concentrations

Considering now just problem ③ (i.e. just  $p(x)$ ):

$$\frac{\partial}{\partial x} (\mu_p^{eff} p E - D_p \frac{\partial p}{\partial x}) + \alpha p = g \quad (p = p^{(k+1)}(x) \text{ unknown}, E = -\frac{\partial \psi}{\partial x} \text{ known})$$

$$-\mu_p^{eff} g_p e^{-\frac{\psi}{V_{th}}} \frac{\partial \psi}{\partial x} - \mu_p^{eff} V_{th} \left( \frac{\partial g_p}{\partial x} e^{-\frac{\psi}{V_{th}}} - \frac{\partial \psi}{\partial x} \frac{1}{V_{th}} g_p e^{-\frac{\psi}{V_{th}}} \right)$$

$$\implies \frac{\partial}{\partial x} \left( -D_p e^{-\frac{\psi}{V_{th}}} \frac{\partial g_p}{\partial x} \right) + \alpha e^{-\frac{\psi}{V_{th}}} g_p = g$$

$$\frac{\partial}{\partial x} \left( -D_p' \frac{\partial g_p}{\partial x} \right) + \alpha' g_p = g$$

Given the **positiveness** of all quantities involved, maximum principle can now be applied to the problem:

$$\exists! g_p \in H^1(\Omega) \text{ s.t. } g_p > 0$$

hence  $p$  (and  $n$ ) is unique and strictly positive

Note how with this notation change (which goes under the name of **Cole-Hopf transformation**) we changed a drift-diffusion current into a purely diffusive current:

$$J_p = q \mu_p^{eff} p E - q D_p \frac{\partial p}{\partial x} = -q D_p' \frac{\partial g_p}{\partial x}$$

Also the opposite can be done. Noting that:

$$\frac{\partial p}{\partial x} = \frac{\partial}{\partial x} \left[ n_i e^{\frac{\phi_p - \psi}{V_{th}}} \right] = \frac{n_i}{V_{th}} e^{\frac{\phi_p - \psi}{V_{th}}} \left( \frac{\partial \phi_p}{\partial x} - \frac{\partial \psi}{\partial x} \right) = \frac{p}{V_{th}} \left( \frac{\partial \phi_p}{\partial x} + E \right)$$

substituting in the drift-diffusion equation we get:

$$J_p = q \mu_p^{eff} p E - q D_p \frac{\partial p}{\partial x} = q \mu_p^{eff} p E - q \mu_p^{eff} V_{th} \frac{p}{V_{th}} \left( \frac{\partial \phi_p}{\partial x} + E \right) = -q \mu_p^{eff} p \frac{\partial \phi_p}{\partial x} = q \mu_p^{eff} p E_p = \text{effective electric field}$$

thus changing a drift-diffusion current into a purely conductive current

Finally, note that:  $g_p \xrightarrow{\psi \rightarrow 0} p$  and  $E_p \xrightarrow{p \rightarrow n_i} E$



Finite element discretization of the current continuity equation

$$\int_0^L \frac{\partial J_p}{\partial x} \phi + \int_0^L x p \phi = \int_0^L g \phi \quad \phi \in H^1(\Omega) \quad x \in \Omega = (0, L)$$

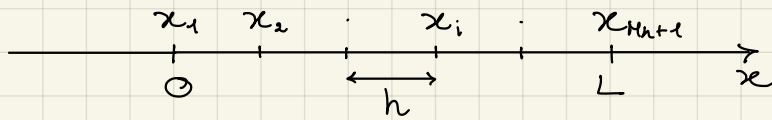
$$\int_{\Omega} \phi \vec{J}_p \cdot \vec{n} - \int_0^L J_p \frac{\partial \phi}{\partial x} + \int_0^L x p \phi = \int_0^L g \phi$$

(w)  $\phi(0) \cdot \vec{J}_p \cdot \vec{n}|_{x=0} + \phi(L) \cdot \vec{J}_p \cdot \vec{n}|_{x=L} - \int_0^L J_p \frac{\partial \phi}{\partial x} + \int_0^L x p \phi = \int_0^L g \phi$

$$p \in V := H^1(\Omega)$$

space discretization

$$p_h \in V_h \subset V$$



$\forall K \in \mathcal{T}_h: J_{p_h} = \text{const over } K$

Considering now  $i = 3$ :

$$-\int_{x_2}^{x_3} J_p \frac{\partial \phi_3}{\partial x} - \int_{x_3}^{x_4} J_p \frac{\partial \phi_3}{\partial x} + \int_{x_2}^{x_3} x p \phi_3 + \int_{x_3}^{x_4} x p \phi_3 = \int_{x_2}^{x_3} g \phi_3 + \int_{x_3}^{x_4} g \phi_3$$

$$-\int_{x_2}^{x_3} J_{p_2} \frac{1}{H_2} - \int_{x_3}^{x_4} J_{p_3} \left(-\frac{1}{H_3}\right) + \int_{x_2}^{x_3} x p \phi_3 + \int_{x_3}^{x_4} x p \phi_3 = \int_{x_2}^{x_3} g \phi_3 + \int_{x_3}^{x_4} g \phi_3$$

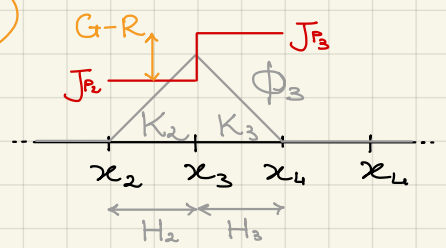
$J_p = J_{p_2}$   
 $\frac{\partial \phi_3}{\partial x} = \frac{1}{H_2}$   
 over  $K_2$

$J_p = J_{p_3}$   
 $\frac{\partial \phi_3}{\partial x} = -\frac{1}{H_3}$   
 over  $K_3$

$$J_{p_3} - J_{p_2} + \int_{x_2}^{x_3} x p \phi_3 + \int_{x_3}^{x_4} x p \phi_3 = \int_{x_2}^{x_3} g \phi_3 + \int_{x_3}^{x_4} g \phi_3$$

$R_3$                        $G_3$

$\implies J_{p_3} - J_{p_2} = G_3 - R_3$



Discrete conservation law

It can be seen as Kirchoff's Current Law:

Hence, the finite element method applied to the current continuity equation embodies, under this assumption, a set of KCLs to be solved for each node of the discretized domain.

$R$  and  $G$  can be computed using a quadrature formula, i.e. a numerical algorithm that approximates the exact integral of a function:

$$I(f) = \int_a^b f(x) dx \sim Q(f)$$

For example,  $Q$  can be the trapezoidal rule:

$$Q(f) := (f(a) + f(b)) \frac{b-a}{2}$$

The advantage of using this formula for the computation of  $R$  and  $G$  comes from the simple expression of  $\phi$ :

$$\int_{x_2}^{x_3} f \phi_3 + \int_{x_3}^{x_4} f \phi_3 \approx (f(x_2) \phi_3(x_2) + f(x_3) \phi_3(x_3)) \frac{H_2}{2} + (f(x_3) \phi_3(x_3) + f(x_4) \phi_3(x_4)) \frac{H_3}{2} \\ = f(x_3) \frac{H_2 + H_3}{2}$$

Hence the conservation law of the previous case becomes:

$$J_{P_3} - J_{P_2} + x_3 p_3 \frac{H_2 + H_3}{2} = g_3 \frac{H_2 + H_3}{2}$$

As a last step, we need an expression for the constant current density  $J_{P_i}$ :

$$J_{P_i} = -q \cdot \frac{D_p}{H_i} \cdot \left[ \text{Be} \left( \frac{\psi(x_{i+1}) - \psi(x_i)}{V_{th}} \right) p_{i+1} - \text{Be} \left( - \frac{\psi(x_{i+1}) - \psi(x_i)}{V_{th}} \right) p_i \right]$$

Scharfetter-Gummel piecewise constant current model

It can be noted that the conservation law equation eventually depends on  $p_{i-1}$ ,  $p_i$  and  $p_{i+1}$  only (as one would naturally expect). This means that the finite element matrix  $\underline{K}_p$ , whose rows represent the conservation law of each node, must necessarily be a tridiagonal matrix.

Def. (M-matrix)

$$\left. \begin{array}{l} \underline{A} \in \mathbb{R}^{n \times n} \text{ invertible w/ } \underline{A}^{-1} \geq 0 \\ (\underline{A})_{ii} > 0 \\ (\underline{A})_{ij} \leq 0 \quad (i \neq j) \end{array} \right\} \implies \underline{A} \text{ is an M-matrix}$$

"Monotone"  
↑



## Theorem

Assume that  $\underline{A} \in \mathbb{R}^{n \times n}$  satisfies the following properties:

$$(a) (\underline{A})_{ii} > 0$$

$$(b) (\underline{A})_{ij} \leq 0 \quad (i \neq j)$$

$$(c) \sum_{i=1}^n (\underline{A})_{ij} \geq 0 \quad \forall j$$

$$(d) \exists j^* \text{ s.t. } \sum_{i=1}^n (\underline{A})_{ij^*} > 0$$

Then  $\underline{A}$  is an M-matrix and  $\underline{A}^{-1} > 0$ .

It is possible to demonstrate that  $\underline{K}_p$  is an M-matrix and satisfies the theorem above.

Being the load vector  $g$  (given by the generation term) strictly positive, we can conclude that the linear system:

$$\underline{K}_p p = g$$

returns only positive values of the concentration  $p$ . So the concentrations ( $p$  and  $n$ ) at every iteration step will be positive and well-defined.

Some concluding comments:

- Gummel's map uses Boltzmann's statistics (a) and (b) to precondition the result of the problem, thus allowing the use of a decoupled (i.e. sequential, without any system of equations) method, whereas Newton's method was fully coupled.
- Gummel's method is very well conditioned, thanks to the fact that the relevant variables ( $\varphi_n$ ,  $\varphi_p$  and  $\varphi$ ) are of the same type and magnitude.
- The advantage of this method, with respect to Newton's, is that it is very robust, does not need any particular requirement to work and it converges for (almost) any initial guess (i.e. it benefits from global convergence). The disadvantage is that its convergence is not as fast.
- It can be demonstrated that Gummel's method becomes much slower for larger devices (i.e. bigger domain  $\Omega$ ).

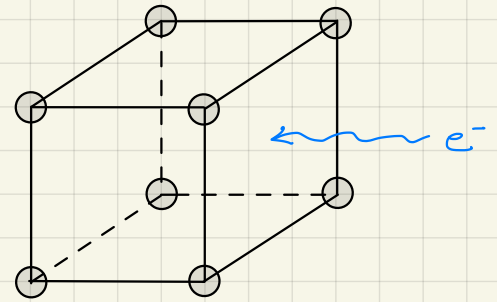
## Mobility models

As we have already discussed, we so far assumed electrical mobility to be constant while in truth it is not.

$$\vec{v}_d = \mu^{el} \vec{E} = \frac{q\tau}{m^*} \vec{E}$$

$m^*$ : effective mass of the particle

$\tau$ : average scattering time

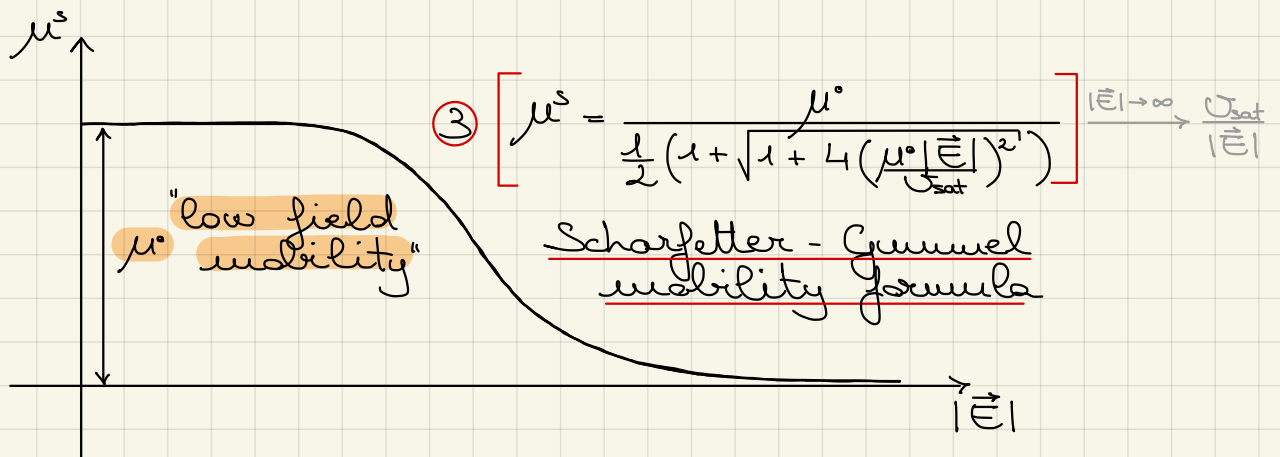


$\tau$  varies with temperature  $T$  (since atoms in the lattice vibrate thus affecting the collision rate and intensity) and with the number of impurities  $N_D^+$  and  $N_A^-$  (since collisions become more frequent).

Furthermore, it would seem that  $|\vec{v}_d|$  diverges with  $|\vec{E}|$ ; however in reality it saturates at a certain saturation velocity  $v_{sat}$  for high electric fields.

① 
$$\mu^L = \mu^0 \left( \frac{T}{T_{ref}} \right)^{-\kappa}$$

② 
$$\mu^{LI} = \mu_{min} + \frac{\mu^L - \mu_{min}}{1 + \left( \frac{N_D^+ + N_A^-}{N_{ref}} \right)^\beta}$$



$$\mu^{LIS} = \frac{\mu^{LI}}{\frac{1}{2} \left( 1 + \sqrt{1 + 4 \left( \frac{\mu^0 |\vec{E}|}{v_{sat}} \right)^2} \right)}$$

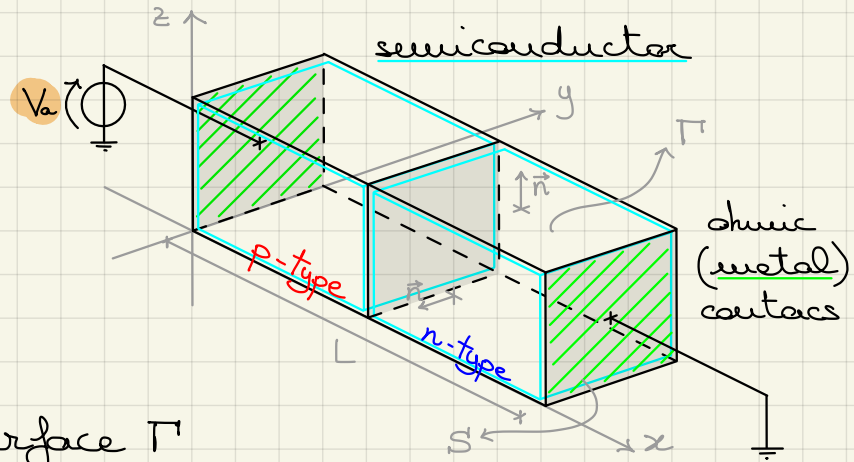
Remember that these mobilities may also vary with  $x$ !

# Semiconductor devices

## P-N JUNCTION

Confinement of electrical phenomena:

$$\left. \begin{aligned} \vec{J}_n \cdot \vec{n} &= 0 \\ \vec{J}_p \cdot \vec{n} &= 0 \\ \vec{D} \cdot \vec{n} &= 0 \end{aligned} \right\} \text{ on lateral surface } \Gamma$$



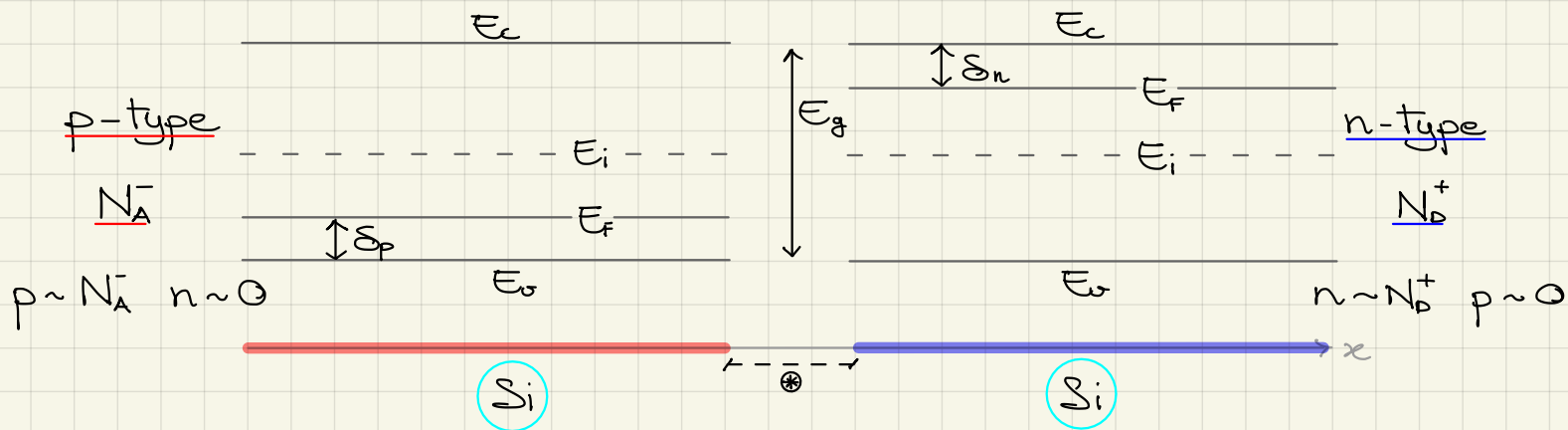
$V_a = 0$ : thermal equilibrium

$V_a > 0$ : forward bias

$V_a < 0$ : reverse bias

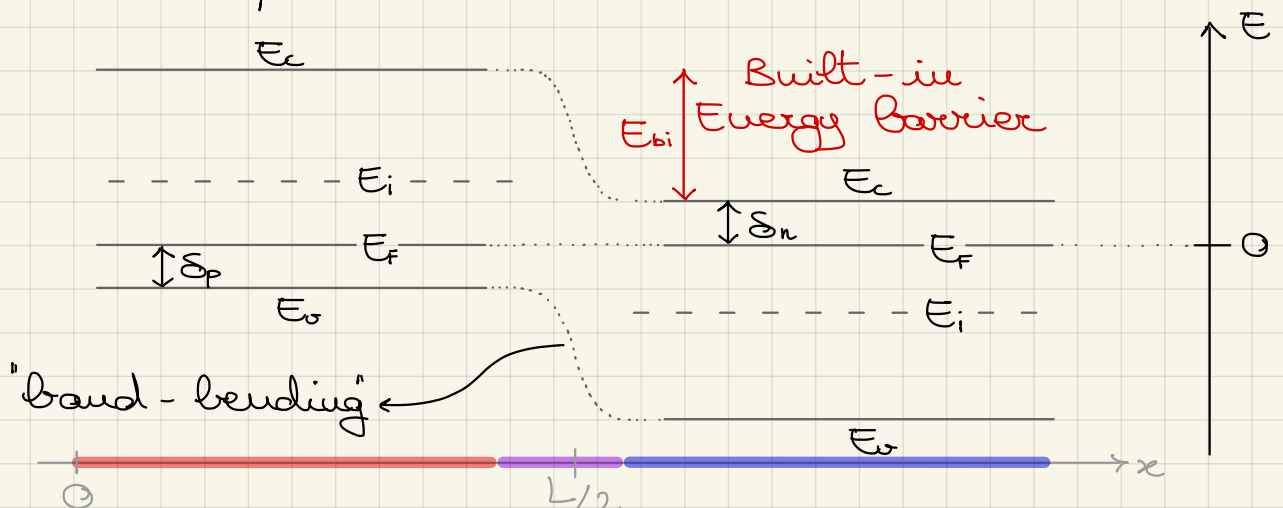
Doping and external voltage affect the energy band diagram.

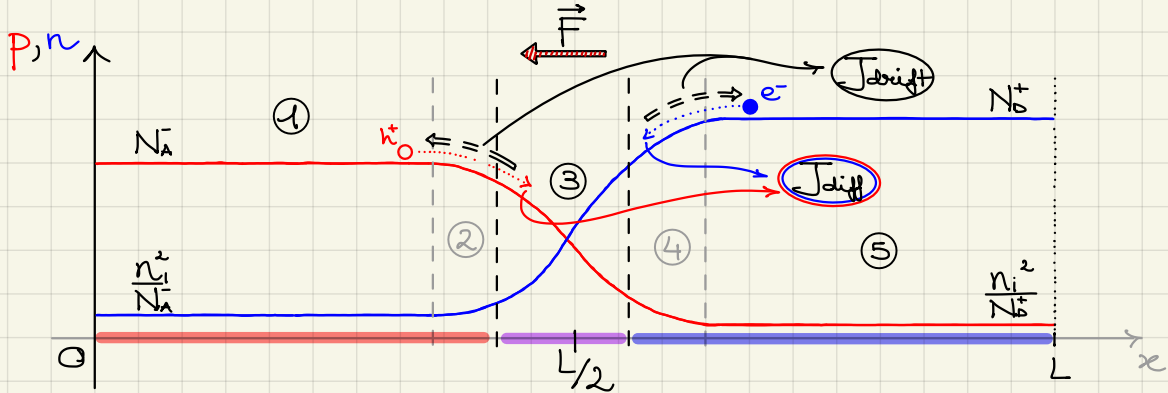
With no voltage applied, the two separated<sup>⊕</sup> regions have their individual band diagram:



When forming a p-n junction, the two diagrams have to share the same Fermi level, which is determined by  $V_a$ .

## Thermal equilibrium





$$p(x) = n_i e^{\frac{E_i(x) - E_F}{k_B T}}$$

$$n(x) = n_i e^{\frac{E_F - E_i(x)}{k_B T}}$$

①, ⑤: (quasi-)neutral regions  $\rightarrow \rho = q(p - n + N_D^+ - N_A^-) \approx q(-N_D^+ + N_D^+) \approx 0$

③: depleted region

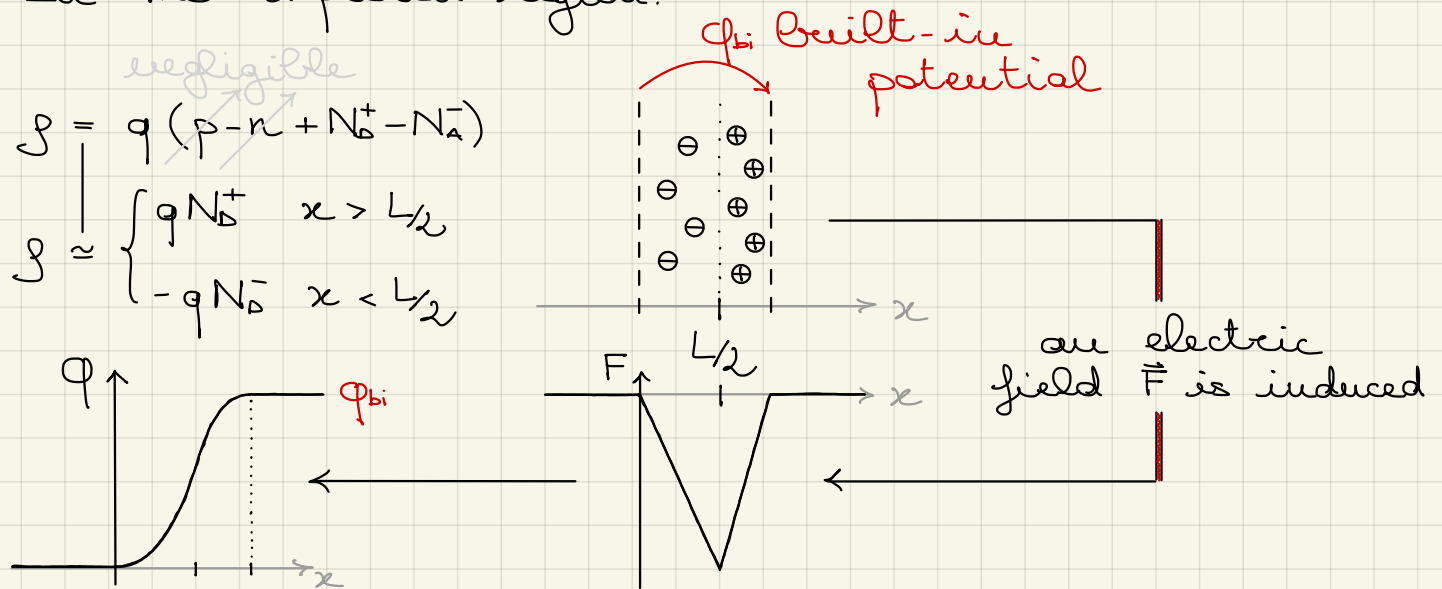
②, ④: minority diffusing regions  $\rightarrow$  typically merged with ① and ⑤ because of complex physical description and marginal relevance

In the depleted region:

$$\rho = q(p - n + N_D^+ - N_A^-)$$

*negligible*

$$\rho \approx \begin{cases} qN_D^+ & x > L/2 \\ -qN_D^- & x < L/2 \end{cases}$$



The built-in potential and energy barrier are strictly related.

$$\begin{aligned} E_{bi} &= E_i|_{x=0} - E_i|_{x=L} = (E_i|_{x=0} - E_F) + (E_F - E_i|_{x=L}) = \\ &= k_B T \ln\left(\frac{N_A^-}{n_i}\right) + k_B T \ln\left(\frac{N_D^+}{n_i}\right) = \underline{k_B T \ln\left(\frac{N_A^- N_D^+}{n_i^2}\right)} \end{aligned}$$

$$\underline{\Phi_{bi}} = \frac{E_{bi}}{q} = \frac{k_B T}{q} \ln\left(\frac{N_A^- N_D^+}{n_i^2}\right) = \underline{V_{th} \ln\left(\frac{N_A^- N_D^+}{n_i^2}\right)}$$

The built-in potential barrier generates a drift current that goes against the natural diffusion current due to the concentration gradients between the n-type and p-type regions.

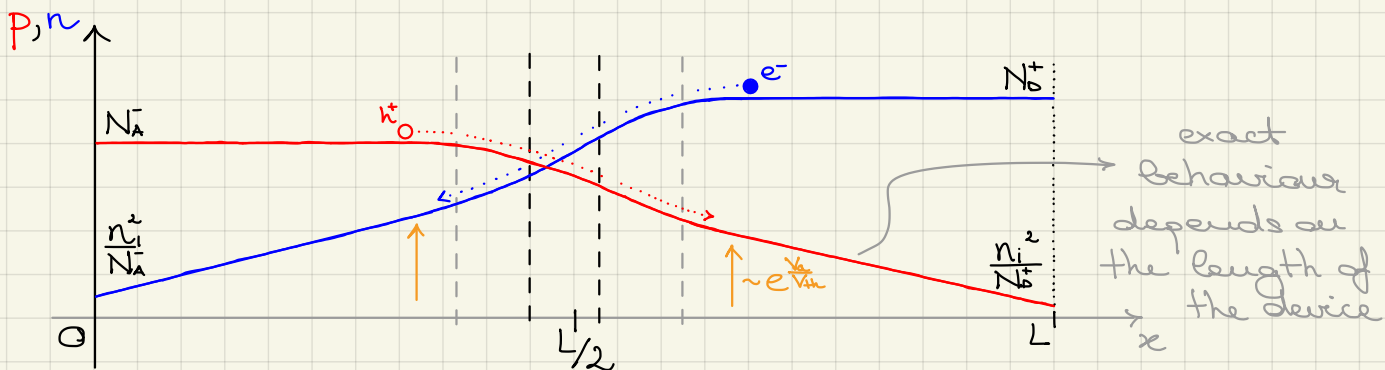
At thermal equilibrium, as one would expect, the TOTAL current is 0, which means that drift and diffusion currents perfectly compensate each other.

(Be careful that in numerical simulations the net current will be therefore given by the subtraction of two large contributions, which might sometimes yield a small but non-zero value).

### • Forward Bias

A forward bias has a direction of the electric field such that it opposes the built-in potential. The barrier is thus reduced.

For  $V_a > \phi_{bi}$ , there is no drift current opposing the diffusion of carriers: holes and electrons can now easily diffuse in the neutral regions.

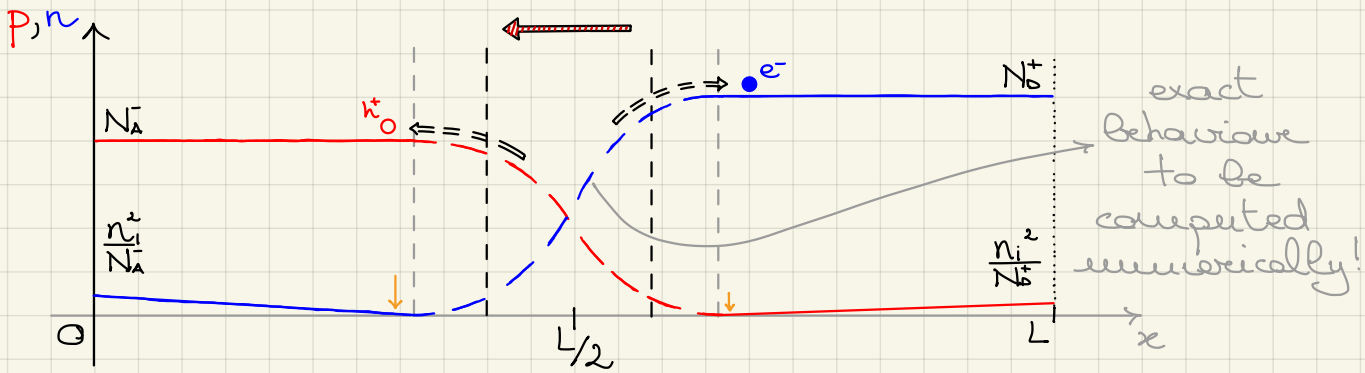


Concentration near the depleted region increases exponentially with  $V_a$ . Carriers diffusing further in the neutral regions recombine with the majority carriers, eventually reaching their equilibrium values.

A diffusion current evidently arises in the neutral region, whose dependency on  $V_a$  has to be exponential.

### • Reverse Bias

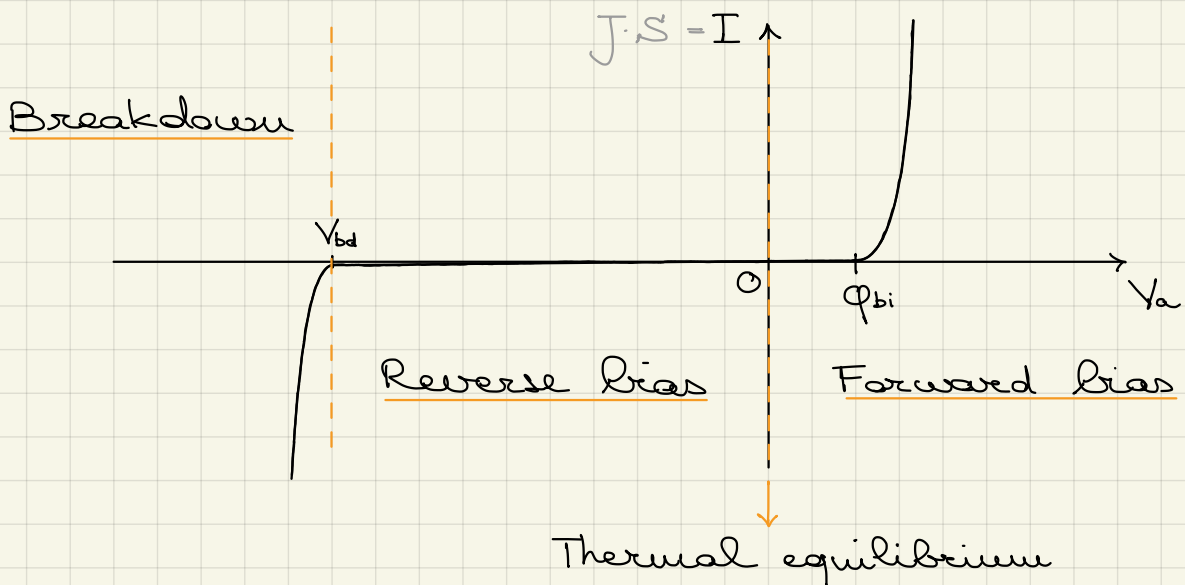
In a reverse bias, the built-in barrier is enhanced by the external voltage. Drift forces in the depleted region are now much stronger than diffusive trends. Almost any carrier reaching depleted region is dragged against the gradient: holes and electrons are basically confined in the p-type and n-type regions, respectively.



Concentration of minority carriers near the depleted region is almost nil since they are dragged by the intense electric field.

A small diffusion current is present in the neutral regions, which is almost constant with respect to  $V_a$ .

(An additional mechanism called breakdown may occur in reverse bias when  $V_a \ll 0$ . The electric field in the depleted region becomes so strong that an accelerated electron (or hole) may carry enough energy to free another electron when colliding with an atom of the lattice. The freed electron, together with the corresponding hole, produces the same effect as other electrons, thus triggering an avalanche mechanism that produces a large amount of current.)



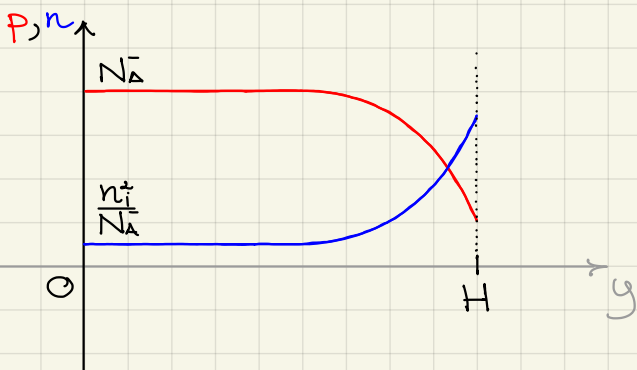
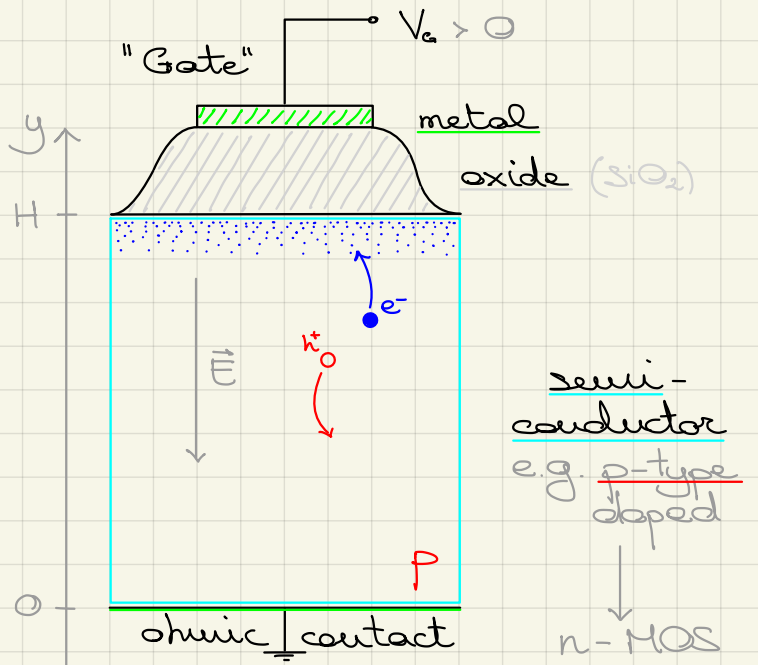


• MOS CAPACITOR

p-type doped substrate:  
n-MOS

n-type doped substrate:  
p-MOS

A variation of  $V_g$  attracts or repels free carriers in the semiconductor.

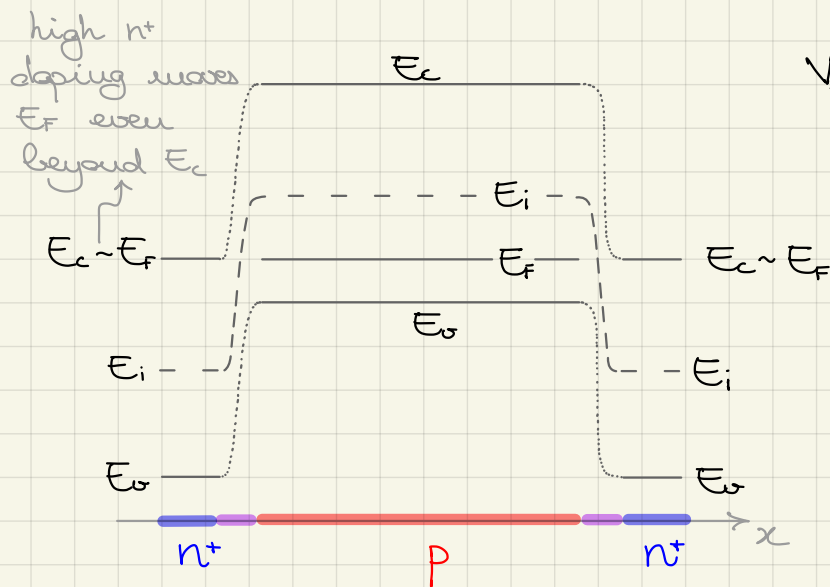
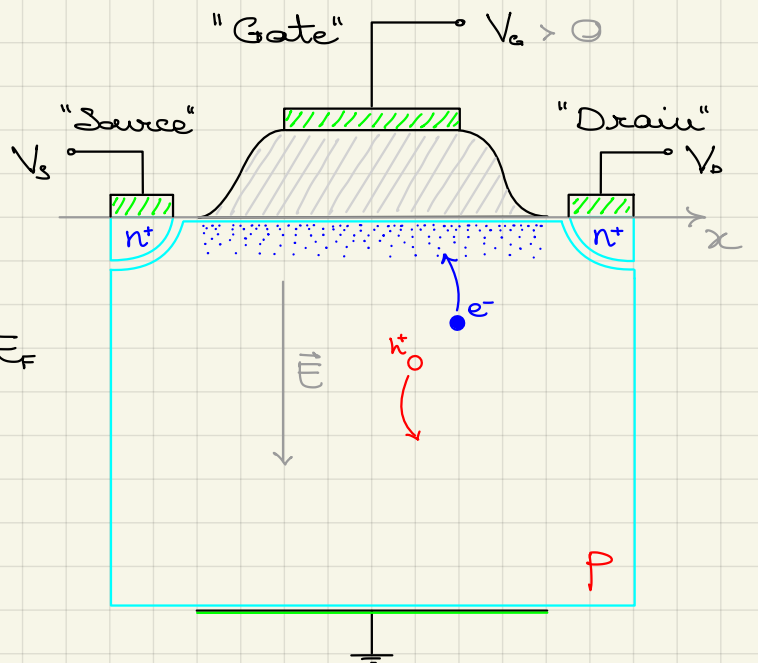


In a n-MOS,  $V_g > 0$  repels holes from the oxide interface and attracts electrons.

For a sufficiently high voltage called **threshold voltage  $V_T$** , electron concentration becomes equal to that of dopants.

A electron channel is formed. By applying a transversal electric field from side to side of the device, current is allowed to flow. The higher the gate voltage, the higher the electron concentration in the channel, the stronger the current for the same applied transversal voltage.

→ • MOS TRANSISTOR





Increasing  $V_g$  will reduce the barrier from source to drain.

Increasing  $|V_D - V_S| = |V_{DS}|$  will help electrons move from source to drain.

The device effectively works as a modulated resistor between source and drain, whose value is determined by the gate.

Note: boundary conditions of ohmic contacts

metal contact + equilibrium & electro-neutrality

$$\left\{ \bar{\varphi}_n = \bar{\varphi}_p = V_{ext} \right\} \quad \left\{ \bar{p} \cdot \bar{n} = n_i^2 \quad \bar{p} - \bar{n} + N_D^+ - N_A^- = 0 \right\}$$

$$\left\{ \bar{\psi} = \bar{\varphi}_n + V_{th} \ln\left(\frac{\bar{n}}{n_i}\right) = \bar{\varphi}_p - V_{th} \ln\left(\frac{\bar{p}}{n_i}\right) \right\}$$