



Appunti universitari

Tesi di laurea

Cartoleria e cancelleria

Stampa file e fotocopie

Print on demand

Rilegature

NUMERO: 2421A

ANNO: 2019

A P P U N T I

STUDENTE: Ricci Francesco

MATERIA: Synthesis and optimization of digital systems - Prof.
Calimera Andrea

Il presente lavoro nasce dall'impegno dell'autore ed è distribuito in accordo con il Centro Appunti.

Tutti i diritti sono riservati. È vietata qualsiasi riproduzione, copia totale o parziale, dei contenuti inseriti nel presente volume, ivi inclusa la memorizzazione, rielaborazione, diffusione o distribuzione dei contenuti stessi mediante qualunque supporto magnetico o cartaceo, piattaforma tecnologica o rete telematica, senza previa autorizzazione scritta dell'autore.

**ATTENZIONE: QUESTI APPUNTI SONO FATTI DA STUDENTIE NON SONO STATI VISIONATI DAL DOCENTE.
IL NOME DEL PROFESSORE, SERVE SOLO PER IDENTIFICARE IL CORSO.**

Introduction

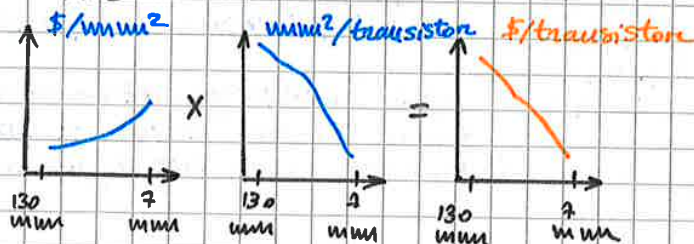
System on chip are everywhere, since electronics brings good value to a lot of domains and applications for a small or moderate cost.

Most of consumer electronics like PC and smartphones are not anymore growing in sales \Rightarrow market saturated... The industry is moved to a **new market driver: IoT!**

The two domains are very different: ICs for desktop are large, power hungry and expensive, while IoT devices must embed more functions (sensors + RF) in a smaller footprint for a much lower price (and usually must be energy efficient as well: think about wearables or biomic devices).

How has been possible to produce more complex and smaller devices for such a lower price? There are two main **enabling factors**:

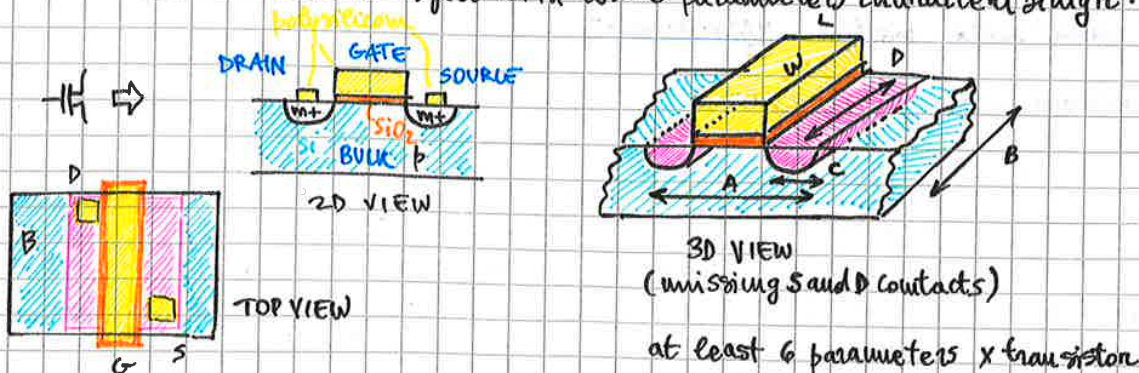
- 1) **Technology Scaling**: silicon vendors are able to provide more complex silicon devices with a larger number of transistors in a smaller area \Rightarrow **Moore's Law**: at each new technology node (130 nm, 90 nm, 65 nm, 45 nm, 32 nm, 14 nm, 7 nm) the overall **cost/mm²** is increasing, but since the geometries of transistors are shrinking, you can pack in the same amount of area more transistors (hence **mm²/transistor** decreases); as a result effect the **cost/transistor** is reducing exponentially over time



So now more complex circuits can be implemented on a smaller footprint. But more complex means **more design complex!** Bigger complexity require more time for designing, increasing the time to market. Today, missing a sweet spot in the market due to delays leads to huge economic losses! How to maintain time to market small, even if complexity grows?

- 2) **Design automation**: it is not possible anymore to design by hand an electronic device, due to its complexity. But it is not only about design: **verification and testing** can easily take 60% of design time \Rightarrow you take more time to design the testbench wrt time to design the circuit.

A transistor is a 3D object with some parameters characterising it:



$$\text{Total cost} = \text{fixed cost} \left(\begin{array}{l} \text{fabrication setup} \\ \text{masks} \\ \text{design cost} \end{array} \right) + \text{variable costs} \left(\begin{array}{l} \text{testing} \\ \text{raw materials} \\ \text{energy/electricity} \\ \text{maintenance} \end{array} \right)$$

Remember that design automation improves QoR (quality of results) and manages complexity:

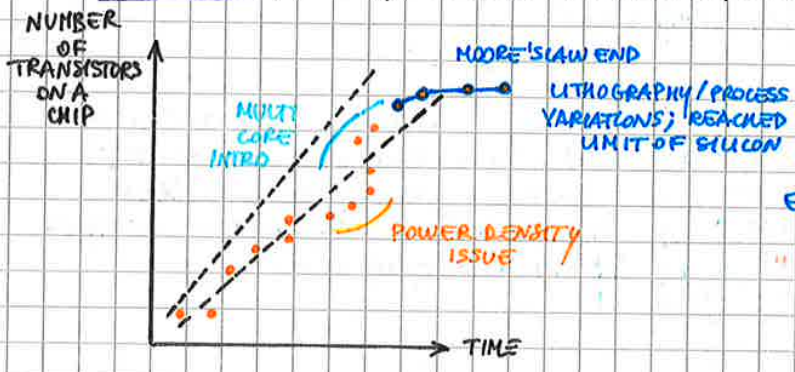
HIGH VOLUME PRODUCTION : Variable cost is the main contribution to final cost : design automation is able to manage complex designs \rightarrow they'll be characterized by higher quality; metrics can be optimized to use the least amount of resources \times IC \rightarrow variable cost decrease

LOW VOLUME PRODUCTION : fixed cost is dominant in final cost : design automation suite reduces the design time, making fixed costs lower

design automation can make the device easier to test, reducing again variable costs

empirical law...
Moore's Law : every 18 months, number of transistors doubles

4/03/18



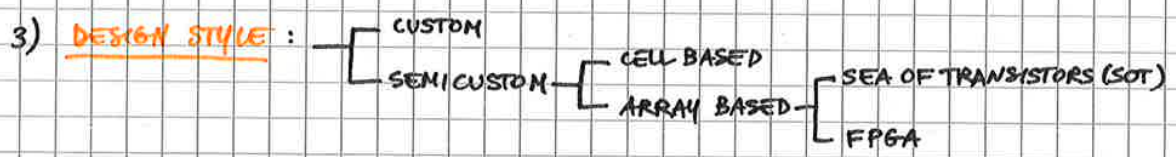
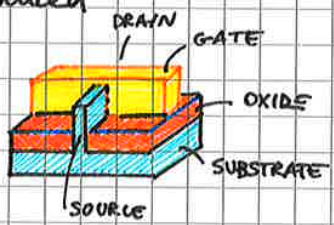
- recap: moving from desktop applications to IoT (# decrease, complexity increase)
- Enabling factors:
- 1) Tech scaling (\$/tr decrease)
 - 2) Design automation (\$/design \downarrow)
 - + manage complex design with less human interaction
 - + improve quality of design (QoD)

c) Silicon on insulator (SOI): bulk is substituted by an insulator, making transistors somehow isolated. This approach mitigates the two issues of noise and leakage.

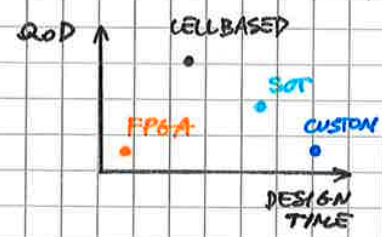


At steady state, a transistor should not have current between source and drain... however a small leakage is present

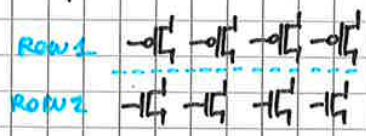
d) thin FAT: gate is a tri-dimensional object, able to better control the transistor: the switch is better, leakage is reduced



a) Custom: define all geometries at layout level: time consuming, complexity, no CAD tables



b) Sea of transistors: a silicon vendor provides a chip with transistors (of the same size) already placed: usually organized in rows of pmos and nmos



at design time you must define all the interconnections, design time reduced

c) FPGA: based on a sea of logic blocks approach (CLBs)



CLB array already designed (LUTs) and placed; programmable interconnections between logic blocks

d) cell-based: the designer can define logic and interconnections; but instead of defining each transistor, you have to use a library of standard cells from a silicon vendor

ALL GATE HAVE THE SAME HEIGHT

For each logic gate multiple layouts are available (the designer is not able to access and see them) each one with a characterization in terms of propagation delay, power consumption, area, etc. → PLACEMENT PROBLEM

good performance, good flexibility

Recap: to design a SoC, we must translate high level specifications into a layout description (that contains all the geometry info to fabricate the device).
As discussed before, design is composed of 3 main phases

1. **Modeling**: describe the system on a chip
2. **Synthesis and optimisation**: provide an implementation
3. **Validation**: check if the implementation corresponds to the model

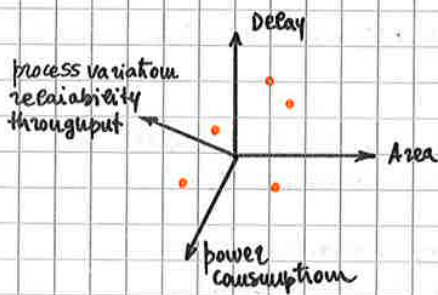
REPEATED AS MANY TIMES AS THE NUMBER OF ABSTRACTION LEVELS

SYNTHESIS WITH NO OPTIMIZATION HAS NO VALUE: optimisation is the means to outperform manual design - Optimisation has multiple objectives:

1. Performance (frequency, latency, throughput)
2. Power consumption
3. Area
4. Testability, dependability

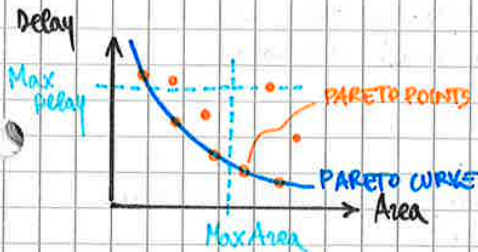
for this reason, it is always a tradeoff: if you want to increase performance, you may want to replicate some functional blocks, increasing area and power consumption, etc.

Design space:



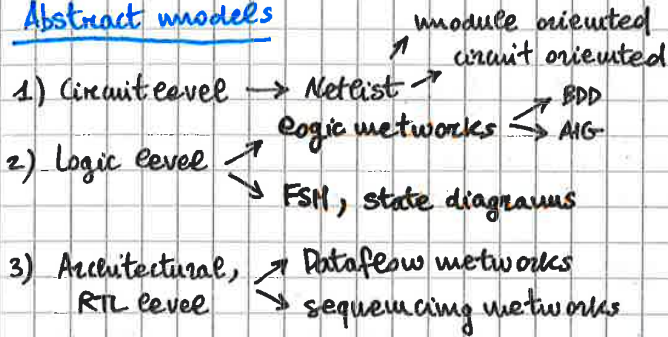
each point represent a possible implementation - How to filter those implementation? A best-all-around solution exist?

Unfortunately, as said before, a solution able to minimize all metrics does not exist. However, given a set of solutions, it is possible to identify some **pareto points** and to trace a **pareto curve** that connects all pareto points.



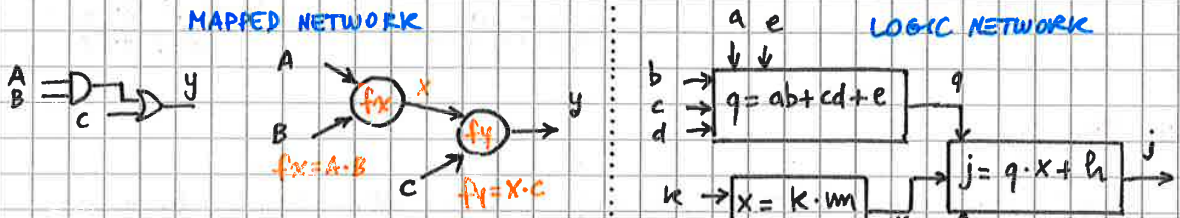
Since a global minimum does not exist, and since an algorithm is not capable to handle all the possible metrics, usually to find/explore pareto solutions, some fixed values are assigned to each metric.

Abstract models



Logic networks:

Mixed of behavioral and structural views, are an interconnection of **nodes** (implementing a Boolean function) and **edges** (representing data dependencies). As always, the model must be acyclic and memoryless. A special case is a **Mapped network**, when nodes correspond to library elements.

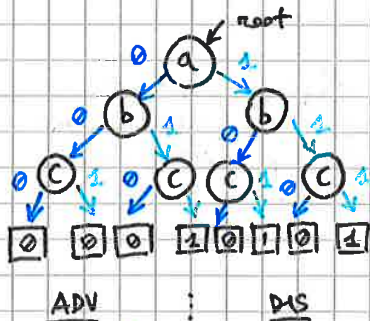


nodes are mapped to gates in the library \Rightarrow one-to-one correspondence between a local function and a gate. (STRUCTURAL VIEW)

Edges represent signal dependency: no info about topology (BEHAVIORAL VIEW)

BDD: Binary Decision Diagrams:

This model is a directed acyclic graph (DAG):
 GRAPH: set of nodes connected by edges
 DIRECTED: edges have direction
 ACYCLIC: no path lead to cycles.



- 1) Each node = decision on a single variable (a, b, or c)
- 2) Value of $f(a,b,c)$ found at leaves (00000...)
- 3) Each path from root to leaf is a row in truth table
- 4) Variable must appear in the same order along each path from root to leaf: by modifying order, the tree may be reduced/optimized (variable reorder).

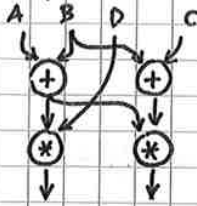
REDUCTION RULES: merge equivalent subtrees
 remove nodes with identical children

- | | |
|--|--|
| <p><u>ADV</u></p> <ol style="list-style-type: none"> 1. Canonical form for each function: only one tree exist. 2. Efficient manipulation of Boolean f. | <p><u>DIS</u></p> <ol style="list-style-type: none"> 1. Size: a tree requires 2^m nodes, with m number of inputs <p>EXPONENTIAL GROWTH!</p> |
|--|--|

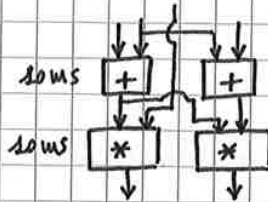
Architectural Level Synthesis CH 4

As already said, the higher the abstraction level the **faster the design**. However, selecting an high abs. level allows to reduce specification of details in the requirements, as well as simplifying the documentation production for a given design. High level description allow also easy modifications and extensions. Those advantages comes with a tradeoff: the **design space is much larger** and synthesis tool are more difficult and expensive to develop.

Example:

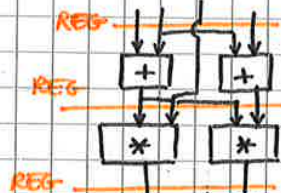


FIRST IMPL.



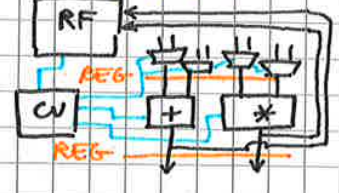
TAKE = 20MS
1 CYCLE

SECOND IMPL.



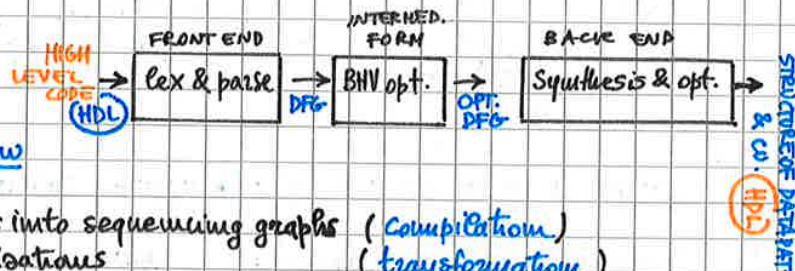
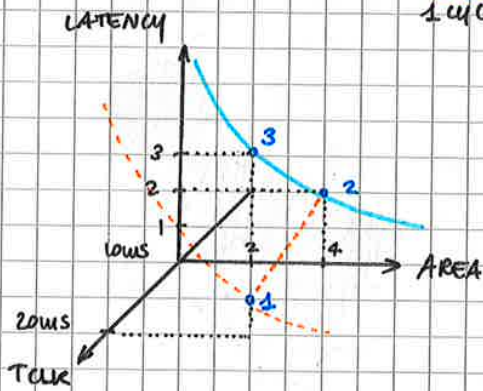
TAKE = 10MS
2 CYCLES

THIRD IMPL.



TAKE = 10MS
3 CYCLES

RESOURCE SHARING BUT NEED CU AND RF



Architectural Level Synthesis Flow

1. Translate HDL models into sequencing graphs (Compilation)
2. Behavioral level optimisations (transformation)
3. Constraint-driven synthesis & Optimisation

1. Compilation: basically it is composed of lexing and parsing to obtain a first tree implementation.

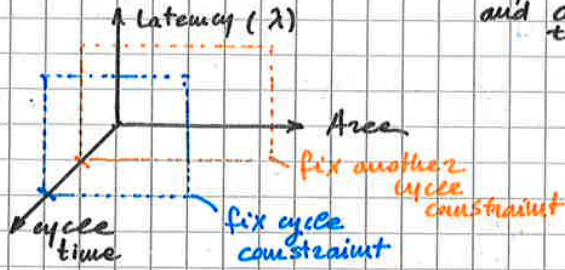
2. Behavioral level optimisations: based on semantic-preserving transformations aiming at simplifying the model. Those transformations can be organised into two categories:

- A) Dataflow based:
- tree-height reduction
 - constant and variable propagation
 - common sub-expression elimination
 - dead-code elimination
 - operator-strength reduction
 - code motion

- B) Control flow based:
- Model expansion
 - conditional expansion
 - Loop expansion
 - Block level expansion

GOAL: explore design space

MULTI-DIMENSIONAL: we consider only Area, latency and cycle (clock period) time



then an algorithm must take in account area and latency

1° APPROACH: fix area constraint and explore latency

2° APPROACH: fix latency constraint and explore Area

We need models for

FIND PARETO POINTS !!

- AREA ESTIMATION:

$$\sum_{i=0}^{\# \text{RESOURCES}} A_i$$

this do not take into account the wiring overhead as well as RF and cv !! Good approximation if RESOURCE DOMINATED IMPLEMENTATION

- LATENCY ESTIMATION

$$\sum_{\text{LONGEST PATH}} d_i$$

delay of each HW resource

this do not take into account delay introduced by the wiring as well as register access and cv delay: good approximation only for RESOURCE DOMINATED IMPL.
NB: $d_0 = 0, d_n = 0$ since root and sink are special nodes useful for formalizing the problem.

ABSOLUTE

ALAP scheduling algorithm : LATENCY CONSTRAINED

This algorithm requires a user threshold on latency: we can provide the minimum latency discovered with ASAP algorithm: $\bar{\lambda} = 3$

USER DEFINED

ALAP ($G(V,E)$, $\bar{\lambda}$) {

schedule v_m by setting $t_m^L = \bar{\lambda} + 1$

repeat {

select a node v_i whose successors are all scheduled

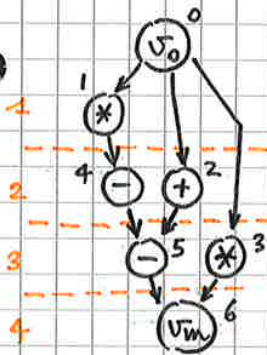
schedule v_i by setting $t_i^L = \min_{j: e_{ij} \in E} \{t_j\} - d_i$

} until (v_0 is scheduled).

return (t^L);

}

Example:



$v_m \Rightarrow t_m = \bar{\lambda} + 1 = 4$

$v_5 \Rightarrow t_5 = \min\{t_m\} - d_5 = 4 - 1 = 3$

$v_3 \Rightarrow t_3 = \min\{t_m\} - d_3 = 4 - 1 = 3$

$v_4 \Rightarrow t_4 = \min\{t_5\} - d_4 = 3 - 1 = 2$

$v_0 \Rightarrow t_0 = \min\{t_1, t_2, t_3\} - d_0 = \min\{4, 2, 3\} - 0 = 1$

REMARKS: ALAP solves a latency constrained problem. latency bound can be set to latency computed by ASAP algorithm: we can define mobility μ

$\mu = t_i^{ALAP} - t_i^{ASAP}$

and it represents the degree of freedom to shift down an operation before violating λ constraint

NB: $\mu_3 = t_3^{ALAP} - t_3^{ASAP} = 3 - 1 = 2$ you can move this node by max 2 positions

$\mu_1 = \mu_4 = \mu_5 = 0$ OPERATION ON LONGEST PATH HAVE ALWAYS MOBILITY = 0

AREA CONSTRAINED:

An area constraint, expressed as the maximum amount of available resources, is present \Rightarrow LATENCY CHANGES DEPENDING ON WHICH KIND OF RESOURCES ARE AVAILABLE, AND IN WHICH QUANTITY

The dataflow graph we want to schedule has the following characteristics:

$$V = \{v_i, i = 0, \dots, m+1\}$$

$$E = \{e_{ij}, i \neq j; i, j = 0, \dots, m\}$$

$$D = \{d_i, i = 0, \dots, m\}$$

$$A = \{a_k, k = 1, \dots, \# \text{RESOURCES}\} \leftarrow \text{for each kind of resource, a maximum number of units is set by the user}$$

Find a labeling function $\varphi: V \rightarrow \mathbb{Z}^+$ such that:

1) $\varphi(v_i) = t_i$

SEQUENCING RELATIONS

2) $t_i \geq t_j + d_j$

k: RESOURCE TYPE
k=1, ..., #RES

MAX NUMB OF RESOURCE OF THAT KIND

RESOURCE BOUNDS

3) $\left| \{v_i \mid T(v_i) = k \text{ and } t_i \leq e \leq t_i + d_i\} \right| \leq a_k$

CARDINALITY

The goal is to minimize the latency, hence minimize t_m (start time of sink)

To solve the area constrained problem, exact and approximated algos are available: let's analyze them one by one.

AREA CONSTRAINED: INTEGER LINEAR PROGRAMMING (ILP)

ILP is a mathematical optimization or feasibility problem: the designer must formalize the scheduling problem as a system of equations.

Some equations represent the constraints, some other identify an objective (or cost) function. Those equations have as unknowns binary decision variables $X = \{x_{i,e}, i = 0, \dots, m, e = 1, \dots, \bar{e}\}$

The advantage of this approach is modularity, since new constraints can be easily added by introducing new equations inside the system, and existing constraints can be modified by rewriting corresponding equations.

The main disadvantage is complexity of design: the designer must write those equations \Rightarrow different problem formulations can lead to different results (hence you may find a local minimum instead of the global one). Moreover, ILP works well up to few thousand variables.

Let's formalize the area constrained scheduling problem using ILP:

1) First, the start time of each operation must be unique

$$\sum_e x_{i,e} = 1 \quad \forall i \in 0 \dots m$$

AREA CONSTRAINED : HU'S ALGO

It is a Greedy strategy-based algorithm, which returns exact solutions under some restrictive assumptions, that are:

- 1) Availability of one single type of resource \rightarrow Co-Processor } MULT
ADD/SUB
DIV
LOGIC
- 2) The HW resource has 1 unit delay: no matter the operations it has to perform, the coprocessor will complete them in a single schedule time cycle.

CO PROCESSORS

HU $(G(V,E), a)$ {

Label the nodes;
 $e = 1$
 Repeat {

$U_e = \{ v_i \mid v_i \text{ has no predecessors or all predecessors have been already scheduled} \}$

$S_e \subseteq U_e$ such that $|S_e| \leq a$ ← # RESOURCES = # CO PROCESSORS

schedule all $v_i \in S_e$

$e++$;

} until $(v_m \text{ is scheduled})$;

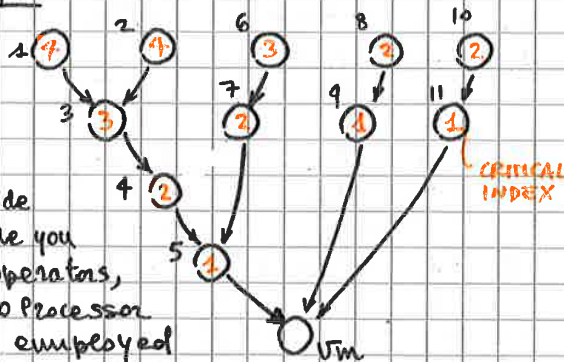
return t_m ;

}

RESOURCES = # CO PROCESSORS

YOU MUST PICK THE MOST CRITICAL OPERATIONS ACCORDING TO A CRITICALITY INDEX CORRESPONDING TO THE DISTANCE FROM SINK.

Example:



NB: inside each node you have no operators, since a Co-Processor must be employed

Suppose #Co-processors = $a = 2$

step	U_e	S_e
$e = 1$	1, 2, 6, 8, 10	1, 2
2	6, 8, 10, 3	3, 6
3	8, 10, 4, 7	8, 10
4	4, 7, 9, 11	4, 7
5	9, 11, 5	9, 11
6	5	5
7	v_m	v_m

$\lambda = 7 - 1 = 6$

Suppose $a = 3$

step	U_e	S_e
1	1, 2, 6, 8, 10	1, 2, 6
2	3, 8, 10, 7	3, 8, 10
3	4, 7, 9, 11	4, 7, 9
4	5, 11	5, 11
5	v_m	v_m

$\lambda = 5 - 1 = 4$

Suppose $a = 1$: then $\lambda = 11$
 Since there are 11 operations of delay = 1 serially scheduled

LATENCY CONSTRAINED : LIST SCHEDULING

minimizes under $\bar{\lambda}$ constraint

The previous list scheduling algorithm can be modified to tackle the complementary problem: given a certain latency, find an implementation that requires the minimum number of resources.

Remember that list scheduling is approximate method, hence the result may not represent the global optimum solution.

↳ LATENCY CONSTRAINT DEFINED BY USER

LIST-R ($G(V,E)$, $\bar{\lambda}$) {

$A = \{1, \dots, 1\}$; initially all resources are set to minimum value = 1.

Run ALAP($G(V,E)$, $\bar{\lambda}$) to compute the latest start time t_i^L for each operation v_i

If ($t_0 < 1$) { return FAIL } : since ALAP was not able to meet the latency constraint $\bar{\lambda}$, there is no chance list scheduling will perform better.

$l = 1$;
repeat {

foreach (resource type k (1, ..., #RESOURCES)) {

READY OPERATIONS → $U_{e,k} = \{v_i \mid v_i \text{ has no predecessors or all predecessors have been executed completely}\}$

$F_{e,k} = \{v_i \mid \text{scheduled before current step and still running}\}$

$S_{e,k} = \{v_i \mid \text{cannot be delayed anymore}\}$

↳ SLACK: $t_i^L - e$

difference between ALAP start time and current step

IF SLACK = 0 v_i MUST BE SCHEDULED, OTHERWISE $\bar{\lambda}$ NOT MET

schedule $S_{e,k}$

If $|S_{e,k}| > a_k$, update a_k

} $e++$;
} until (v_m is scheduled);

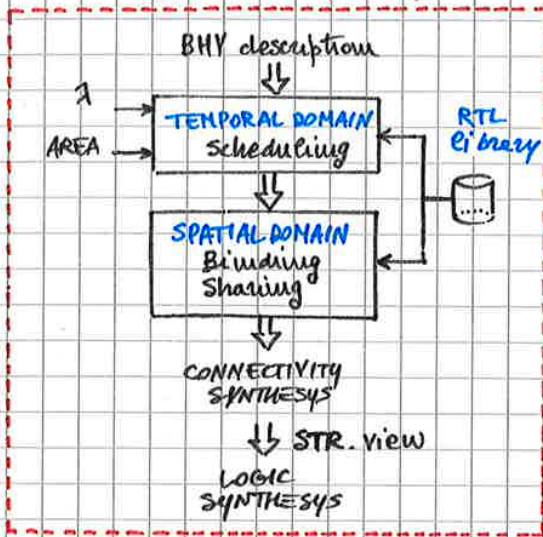
} return (t_m , A)

CH6 Sharing and Binding

In chapter 4 the architectural level synthesis: as already said, the synthesis process can be classified into **TEMPORAL DOMAIN** (scheduling) and **SPATIAL DOMAIN** (binding and sharing).

Scheduling was explained in CH5, let's now focus on sharing and binding

ARCH. LEVEL SYNTHESIS

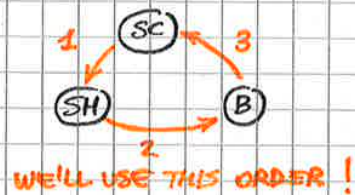
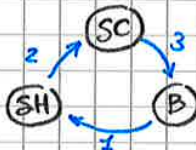
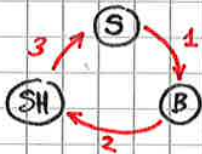


- BINDING: assign a specific HW resources (contained inside the RTL library) to a specific operation in the DFG.

Usually RTL libraries contains multiple implementations/configurations: for example a large & fast or a small and slow. The binding process can then produce different results based on the RTL library content: the best solution can be filtered by looking at some metrics (like AREA, power consumption, etc.) BUT NOT LATENCY (ALREADY DEFINED DURING SCHEDULING).

- SHARING: define two or more operations able to share the same HW resource

Remember the linear dependency between scheduling, binding, sharing



BINDING/SHARING SCENARIOS:

1) **DEDICATED RESOURCES:** minimum λ (of executed ASAP), maximum area (one HW resource per operation)

- ⇒ NO SHARING
- ⇒ BINDING INTRINSICALLY DEFINED

2) **ONE MULTI-TASK RESOURCE:** like an arithmetic coprocessor, able to perform various kind of operations (+, -, <, *, ÷)

The HUI'S ALGO (with area = 1) can be exploited for scheduling

- ⇒ NO BINDING: there is not a direct correspondance between op and HW resource, since there is only one resource that performs everything
- ⇒ MAXIMUM SHARING

3) **MANY MULTI-TASK RESOURCES:** same coprocessor as before, but now more than one is employed; scheduling done with HUI'S ALGO (area > 1).

- ⇒ PARTIAL BINDING
- ⇒ SHARING IS AN OPTIMIZATION PROBLEM (NP complex)

4) **ONE RESOURCE X-TYPE:** no coprocessors, but some address multipliers, divisors, etc. For each component only one configuration is available.

- ⇒ SCHEDULING: LIST SCHEDULING
- ⇒ BINDING: PARTIALLY DEFINED
- ⇒ SHARING: OPTIMIZATION PROBLEM

1. DEFINE PROBLEM: COMPATIBILITY
2. DEFINE ABSTRACT MODEL
3. DEFINE ALGORITHM USING ABSTRACT MODEL

Sharing problem

Starting from a scheduled sequencing graph, where commutivity between operations is well defined, you want to identify two or more operations that can be bound on the same HW resource. The concept of **compatibility** must be introduced.

Two or more operations are compatible (and then can share the same HW) if:

- 1) they can be implemented by the same resource type
- 2) they are not concurrent (not used in the same step \equiv not overlapping)

Clearly, since an operation can be mapped on a certain resource type only, the sharing problem benefits of **problem decomposition**: you can analyse sharing for each resource type. \Rightarrow REDUCE TIME AND COMPLEXITY OF SHARING PROBLEM

Abstract model: how to model compatibility?

COMPATIBILITY GRAPH G^+

COMPLEMENT

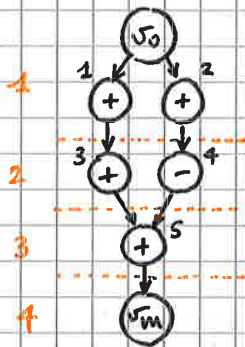
CONFLICT GRAPH G^-

modes: operations in DFG
edges: $e_{ij} \Rightarrow i$ and j are compatible

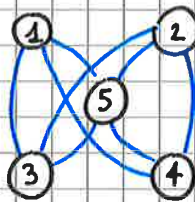


modes: operations in DFG
edges: $e_{ij} \Rightarrow i$ and j are not compatible

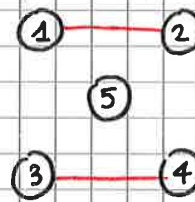
Example:



COMPATIBILITY G^+



CONFLICT G^-



As can be seen, G^+ and G^- are complementary

In order to find the optimum resource sharing you must analyse the G^+ or G^- graphs: in a G^+ graph, it is possible to identify **cliques** (subset of modes that are fully interconnected between each other); maximum sharing correspond to minimum number of cliques. \Rightarrow PARTITIONING

In a G^- graph, a clique is a subset of modes that has no edges between each other; maximum sharing correspond to minimum number of cliques.

Usually this problem is tackled by assigning a color to each mode such that connected modes have different colors. \Rightarrow COLORING

Both partitioning and coloring are NP complex problems. To reduce problem complexity we have to refine the G^+ and G^- graphs into an equivalent representation: a **perfect graph**. Those graphs have polynomial time solution to coloring and partitioning:

- 1) Golumbic's algo (comparability)
- 2) left-edge algo (interval) \leftarrow WE'LL SEE THIS ALGO IN THIS COURSE

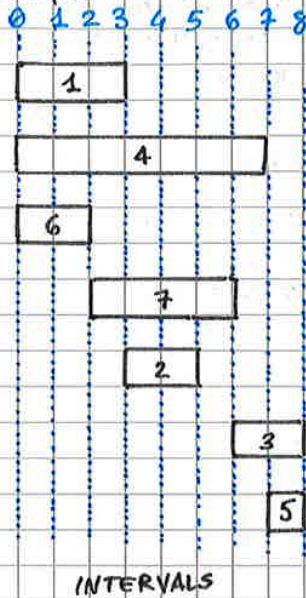
In particular, perfect graphs must specify the "source" of conflict/compatibility by means of **oriented edges** (\rightarrow): a "partial order" is then defined, imposing that the operation at arrow tip must be executed after the one at the other edge end. \Rightarrow THIS ORDER INFO IS DERIVED FROM THE SCHEDULE

Left Edge Algorithm

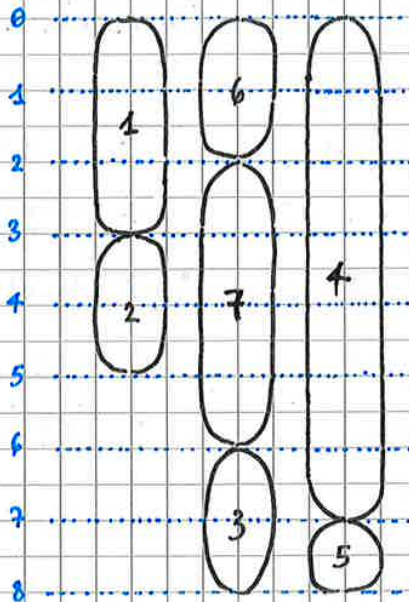
This is a coloring algorithm applied on an interval graph. Hence, the goal is to find the minimum number of colours such that all disjoint nodes share the same colour.

Supported variables $\left\{ \begin{array}{l} L = \{ I_i : \text{intervals ordered w.r.t. left edge} \} \\ C = \{ c_i : \text{set of colours} \rightarrow \text{initially is one colour only} \} \end{array} \right.$

Example:



←
moving from DFG to interval graph



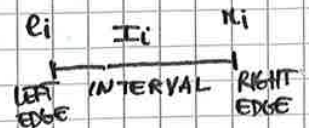
$L = \{ I_1, I_4, I_6, I_7, I_2, I_3, I_5 \}$ NB: is ordered!

Left edge (I)

$\left\{ \begin{array}{l} I: \text{intervals (not sorted)} \\ L: \text{intervals (sorted)} \\ S: \text{set of intervals assigned to a specific colour} \\ r: \text{right edge of the last coloured interval} \\ C: \text{number of colours currently used (initially is 1)} \end{array} \right.$

Supported variables

Sort element of I in a list L in ascending order of e_i ;
 $C = 1$;



while (there are some remaining uncoloured intervals) {

$S = \{ \}$; // empty set

$r = 0$;

while (there is at least one interval I_i such that $e_i \geq r$) {

$S = S \cup I_i$ // add I_i to set S

$L = L - I_i$ // remove I_i from L

$r = r_i$ // update the variable r with the value of right edge of the last coloured interval I_i

}

Assign colour to S

Increment number of colours ($C++$)

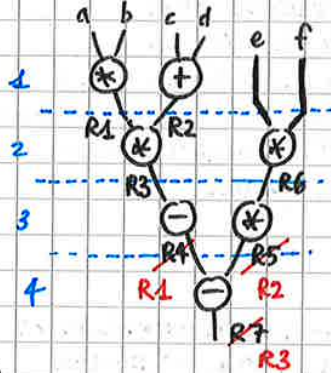
}

SHARE REGS $\left\{ \begin{array}{l} \text{REDUCE CARDINALITY} \\ \text{REDUCE POWER CONSUMPTION} \\ \text{REDUCE SIZE AND COST} \end{array} \right.$

Register Sharing problem

If no sharing is present, it is possible to hardwire the output of a module to an input of the following one (if they are scheduled in adjacent time slots)

However, if a certain degree of sharing is applied, registers becomes crucial! Moreover, registers are expensive and in a complex DFG you may need thousands of them; they switch every clock cycle also, consuming large power! For this reason, a register sharing problem must be considered

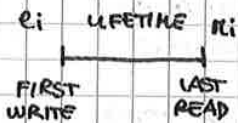


Registers must be used to hold values of variables

It is not possible to share R1 and R3, because the corresponding register has to provide stable inputs for the module (*) but at the same time has to sample its output... Hence to share a register, it must be present a gap of at least 1 time unit (for example R1 and R4).

Let's formalise the problem: register sharing is based on lifetime of variables

Lifetime: interval between the time at which the variable is generated and the latest time at which it is referenced (as input of another operation).



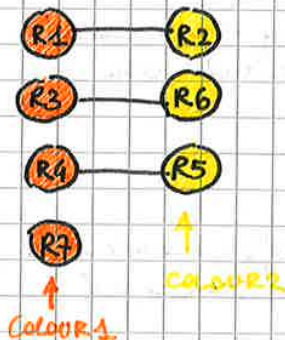
$$\text{Lifetime} = t_{\text{LAST READ}} - t_{\text{PRODUCTION}}$$

This problem can be solved resorting again on perfect graphs: in particular, for an interval graph we must define:

- Nodes \rightarrow variables lifetime intervals
- Edges \rightarrow lifetime overlaps

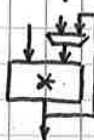
Using left Edge Algo, find minimum number of registers to store all variables

For the circuit above, the interval graph is



Register sharing for iterative bodies

Some architectures have "variable loop": for example in a multiply and accumulate, the result goes as feedback to 1 input of the multiplexer...



To correctly share registers those loop connections must be considered: hence a circular-arc conflict graph is employed.

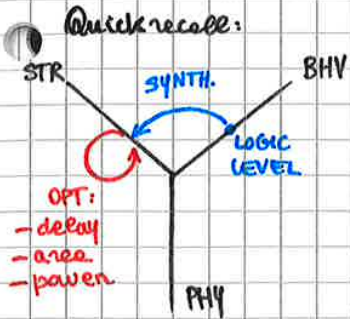
Bus sharing and binding

Find the minimum number of busses to accommodate all data transfer \rightarrow conflict variables need dedicated connection

For the previous example, 1 BUS allows one colour transfer, 2 BUSES allow all variables to be transferred.

NB: after register sharing; you must introduce steering logic (like muxes) to correctly drive the shared regs...

Logic synthesis and optimization CH7



input: behavioural description in HDL
primitives are boolean

output: netlist (structural description) in HDL
↳ connection of logic gates (from tech library)

GOAL: define and explore the design space



Example:

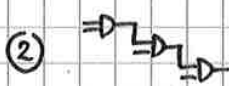
$$f = abcd$$

tech library includes $AND2(2,2)$, $AND3(3,3)$

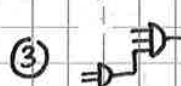
possible implementations:



area = 6
delay = 4



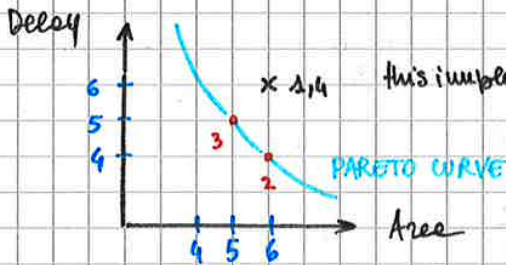
area = 6
delay = 6



area = 5
delay = 5



area = 6
delay = 6



these implementations are not pareto points.

Going forward, in this lecture we'll see some background concepts, some definitions to formulate the problem and finally some exact and approximate algos to solve the problem.

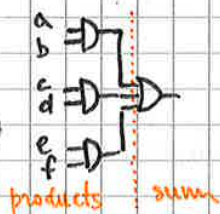
Background:

SYNTHESIS OF TWO LEVEL CIRCUITS

SYNTHESIS OF MULTI LEVEL CIRCUITS

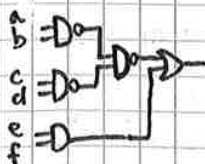
TWO LEVEL: SOP

$$f = ab + cd + ef$$



MULTI LEVEL:

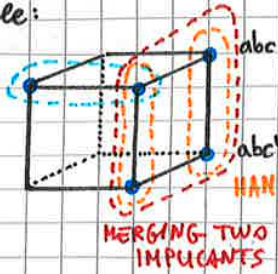
$$f = \overline{ab \cdot cd} + ef$$



no restrictions on kind of gates number of levels

Nowadays no one is using two-level circuits anymore, but two level optimisation is used within multi level implementations

Example:



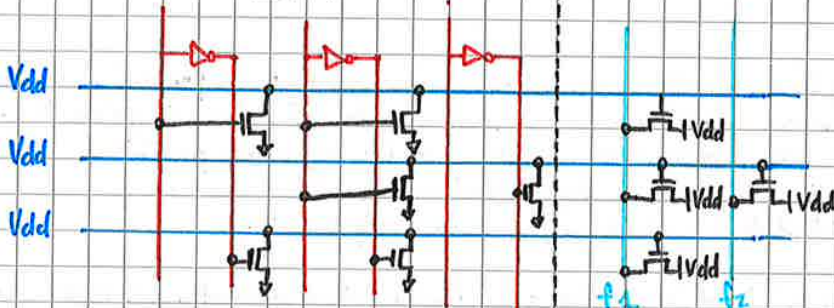
$$abc \cup abc' = abc + abc' = ab$$

PLAs (programmable logic array)

$$f: B^3 \rightarrow B^2 \quad \begin{cases} f_1 = a'b' + b'c + ab \\ f_2 = b'c \end{cases}$$

AND: products

OR: sums



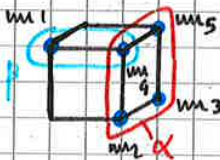
horizontal wires = # products
vertical wires = # literals

horizontal wires = # products
vertical wires = # outputs

If you want to reduce area:

- 1) find a representation with the minimum number of implicants \Rightarrow reduce # products
- 2) within a product, you want to reduce the number of transistors \Rightarrow reduce # literals

Minimum cover: find minimum number of implicants that cover all minterms

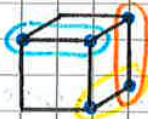


from $f = m_1 + m_2 + m_3 + m_4 + m_5$ to $f = \alpha + \beta$

5 ROWS, 15 TRANS.

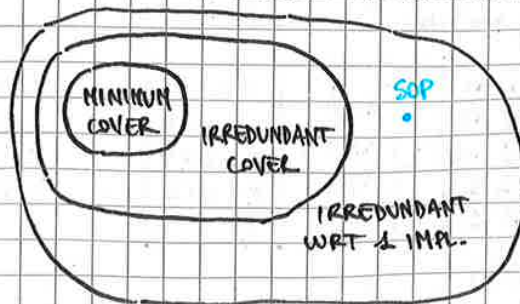
2 ROWS, 3 TRANS

Irredundant cover (minimal cover): none implicant can be dropped



Obviously, if any of the implicant is dropped, at least one minterm is cover anymore. This solution may be not the optimum one.

Irredundant cover with respect of one implicant: none implicant can be contained into another implicant



second step of

The reduction consist on finding an appropriate vector X such that

- $A \cdot X \geq b$
- $e_1(x)$ is minimum $\sum x_i$? ? ? ?

In practice you must select enough columns to cover all constraints, but at the same time you want to minimize the cardinality of X (hence select the minimum number of cols that cover all constraints).

COVERING PROBLEM \Rightarrow NP HARD PROBLEM

in the previous example, $X = [1111]^T$, we obtain $Y = [12221]^T$.

But, by selecting $X = [1101]^T$, we obtain $Y = [12111]^T$.

In both cases, the first requirement ($Y \geq b$) is met, but the second solution has $e_1(x)$ smaller:

$$e_1([1111]^T) = \sum x_i = 1+2+2+2+1 = 8$$

$$e_1([1101]^T) = \sum x_i = 1+2+1+1+1 = 6$$

2-LEVEL OPT.

(slower, high complexity)

EXACT METHODS: they find the minimum cover

Quine McCluskey

Petrick's

Espresso

APPROXIMATE: they find the minimum irredundant cover (weaker but faster, lower complexity)

The exact methods all need the prime table as a starting point; they all reduce the prime table using different approaches:

- 1) Quine McCluskey: tabular methods
- 2) Petrick's: POS clauses
- 3) Espresso: tabular methods + branch & bound

NB: the "reduce" = optimization problem becomes a "covering problem" !

	A	B	C	D	E	F
m1	1	1	0	0	0	0
m4	0	0	0	0	1	0
m5	1	0	0	0	1	0
m6	0	0	0	0	1	1
m7	0	0	0	0	1	1
m9	0	1	1	0	0	0
m11	0	0	1	1	0	0
m14	0	0	0	0	0	1
m15	0	0	0	1	0	1

2° step: all minterms contained in the removed primes are removed from the table

	A	B	C	D
m1	1	1	0	0
m9	0	1	1	0
m11	0	0	1	1

1° step

After the essential primes are removed, we must continue to reduce the table. A rule of thumb suggests to remove the columns with largest number of ones first, because they cover a larger number of minterms. In our example B and C are valuable choices. Let's remove C:

	A	B	C	D
m1	1	1	0	0
m9	0	1	1	0
m11	0	0	1	1

	A	B	D
m1	1	1	0

$$C_{min} = \{A, C, E, F\}$$

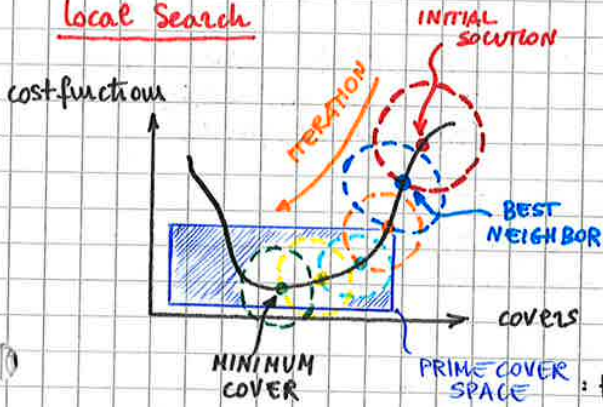
we remove A for example

NB: we found a minimum cover, but it is not the only possible one. For example: $C_{min}^2 = \{B, C, E, F\}$, $C_{min}^3 = \{B, D, E, F\}$, ...

Heuristic optimization

Approximate methods are used as internal engines within multi-level optimization and synthesis tools. For this reason, they must be faster (w.r.t exact methods); hence bottlenecks like prime generation and storage and covering problems are avoided. Even if the heuristic is able to provide "reasonably small" irredundant covers in a "reasonable time", the overall quality of the final cover may be not optimal (local minimum instead of global one).

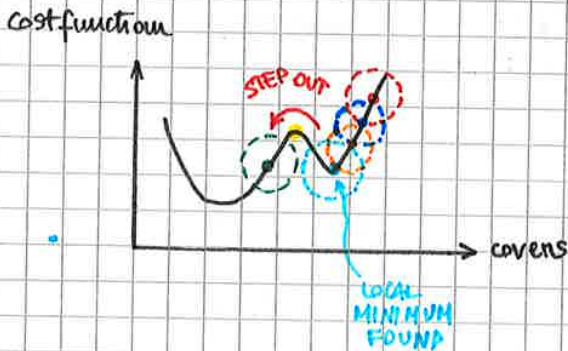
Local Search



Basic strategy:

- 1) Starting point is a random SOP (cover)
- 2) Search for feasible neighbors whose cost is lower
- 3) If one is found, iterate (go back to 1)

PRIME COVER SPACE: the exact methods start here, inside the prime cover space.



Drawback:

- 1) local search may stuck at local minimum (no convex problem)
- 2) Need a mechanism to step out (accept worse solution)

Operators:

- ITERATE
- 1) Expand: make implicants prime → make prime cover
 - 2) Irredundant: make cover irredundant
 - 3) Reduce: reduce size of each implicant while preserving interms covering, returning to a non-prime cover

- Expand and irredundant, make the cover prime and irredundant, therefore the best neighbor is found
- Reduce: perturb the cover finding some worse covers to step out of a local minimum

The possible issues are

- 1) Validity check of expansion
- 2) Find the best "direction" to move to
- 3) Order of expansion (how many literals drop)

Heuristic optimisation tools

- 1) Tautology: check if a function is always true
- 2) Containment: check if a cube is contained in a cover
 this tool is essential to understand if a given cover covers all the minterms. It is also used by the "irredundant" operator to find cubes or implicants contained in other implicants of the cover
- 3) Complementation: compute the complement of a function
 $OFFSET = \text{complement}(ONSET + DCSET)$
 OFFSET is crucial during "expand" operator, to understand if a given implicant is valid

Background:

- ① given a function $f: B^m \rightarrow B^m$, the Shannon decomposition w.r.t. a variable x_i is:

$$f = x_i f_{x_i} + \bar{x}_i f_{\bar{x}_i}$$

↑ cofactor w.r.t. x_i
↑ cofactor w.r.t. \bar{x}_i

where

$$f = f(x_1, x_2, \dots, x_i, \dots, x_m)$$

$$f_{x_i} = f(x_1, x_2, \dots, 1, \dots, x_m)$$

$$f_{\bar{x}_i} = f(x_1, x_2, \dots, 0, \dots, x_m)$$

Example:

$$f = ab + bc + ac = a(bc) + \bar{a}(bc)$$

$$f_a = b + bc + c = b + c$$

$$f_{\bar{a}} = bc$$

- ② a given function f is unate w.r.t. x_i if $f_{x_i} \geq f_{\bar{x}_i}$

if x_i is contained into f always with same polarity, then f is

- x_i always positive in $f \rightarrow$ positive unate
- x_i always negative in $f \rightarrow$ negative unate

Example:

$$f = ab + bc + ac \rightarrow \left. \begin{array}{l} \text{positive unate wrt } a \\ \text{positive unate wrt } b \\ \text{positive unate wrt } c \end{array} \right\} \rightarrow f \text{ is positive unate}$$

$$f = a'b + bc + dc \rightarrow \left. \begin{array}{l} \text{negative unate wrt } a \\ \text{pos unate wrt } b \\ \text{positive unate wrt } c \end{array} \right\} f \text{ is binate}$$

$$f = a'b + bc + ac \rightarrow \left. \begin{array}{l} \text{binate wrt } a \\ \text{pos unate wrt } b \\ \text{pos unate wrt } c \end{array} \right\} f \text{ is binate}$$

Matrix representation of logic covers

Approximate methods do not use prime tables (big size, difficult build); instead a different table approach is used:

the table is composed by one row per implicant and one column per variable.

$$\text{encoding } x_{ij} = \begin{cases} 00 & \text{if intersection between } \text{var}_j \text{ and implicant } i \text{ is void} \\ 01 & \text{if } \text{var}_j \text{ is positive in implicant } i \\ 10 & \text{if } \text{var}_j \text{ is negative in implicant } i \\ 11 & \text{if } \text{var}_j \text{ is a don't care in implicant } i \end{cases}$$

This encoding allows a very easy and easy to build table: one cell contains 4 elements!

Example: $f = a'd' + a'b + ab + ac'd$

	a	b	c	d
a'd'	10	11	11	10
a'b	10	01	11	11
ab	01	01	11	11
ac'd	01	11	10	01

Moreover, binary operations are very fast and can be parallelized:

- BITWISE AND: intersection
- BITWISE OR: supercube
- BITWISE COMPLEMENT: rise a don't care

Example: intersection between $f = a'b'c$ and $g = b'$

	a	b	c
a'b'c	10	10	01
b'	11	10	11
	10	10	01 = a'b'c



intersection between $f = a'b'c$ and $g = b$

	a	b	c
a'b'c	10	10	01
b	11	01	11
	10	00	01



void: means the intersection is empty

Example: don't care rise: from $ac'd$ to $*c'd$

	a	b	c	d
ac'd	01	11	10	01
BITWISE COMPL. OF a	10	00	00	00
	11	11	10	01 = c'd

BITWISE OR
remember:
 $a = 01\ 11\ 11\ 11$
the expected value

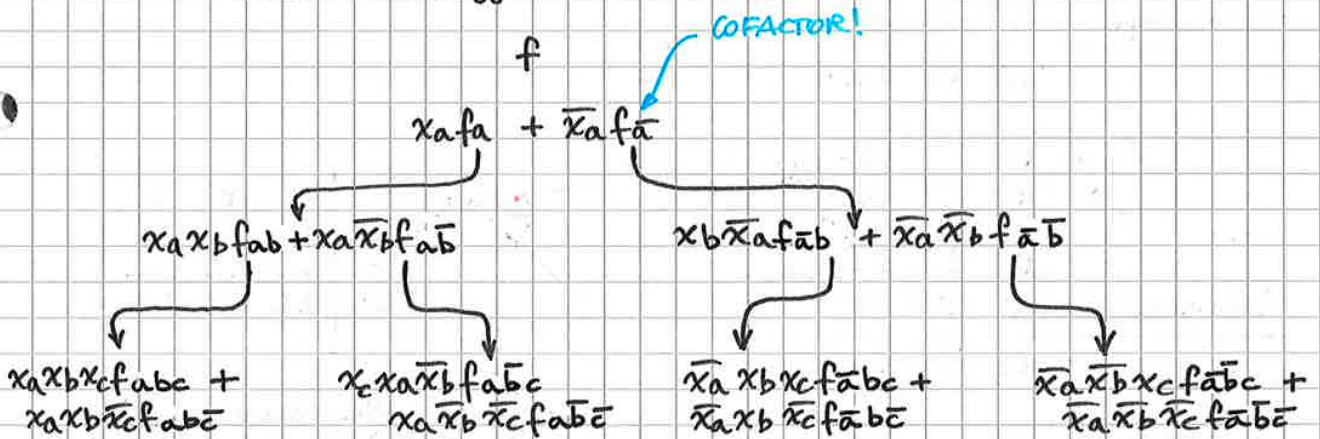
Tautology check

To check if a function is always true, a recursive algorithm must be used. For each variable of the function, the cofactors must be computed: if they all are tautologies, the entire function is a tautology.

Terminial cases

- 1) The cover matrix has a row of ones \Leftrightarrow TAUTOLOGY
- 2) The cover matrix has a column of zeros \Leftrightarrow NOT TAUTOLOGY
- 3) The cover matrix depends on one variable and there is no column of zeros \Leftrightarrow TAUTOLOGY

Example: suppose f depends on 3 variables a, b, c .
To check if f is a tautology, you must explore all the tautology tree:



As soon as one cofactor is not a tautology, the entire function f is not a tautology and therefore no need for computing all remaining cofactors.

If a cofactor is a tautology, you can skip the computation of child cofactors, because they'll be tautologies as well. However, all cofactors must be a tautology to prove that f is a tautology.

Computation of f_{ab}

1° row	11	01	11	AND
b'	11	10	11	
	11	00	11	drop
		VOID		
2° row	11	11	01	AND
b'	11	10	11	
	11	10	01	OR
bitw comple b'	00	01	00	
	11	11	01	
3° row	11	10	10	AND
b'	11	10	11	
	11	10	10	OR
bitw comple b'	00	01	00	
	11	11	10	

$f_{ab} =$

	a	b	c
c	11	11	01
c'	11	11	10

f_{ab} is a tautology, because it depends on a single variable and there is no column of zeros



Since both f_{ab} and f_{ab} are tautologies, f_a is a tautology as well. What is missing is a check on cofactor $f_{a'}$:

Computation of $f_{a'}$

1° row	01	01	11	AND
a'	10	11	11	
	00	01	11	drop
		VOID		
2° row	01	11	01	AND
a'	10	11	11	
	00	11	01	drop
		VOID		
3° row	01	10	10	AND
a'	10	11	11	
	00	10	10	drop
		VOID		
4° row	10	11	11	AND
a'	10	11	11	
	10	11	11	OR
bitw comple a'	01	00	00	
	11	11	11	tautology! (row of ones)

Since f_a and $f_{a'}$ are tautologies, f is a tautology! In fact:



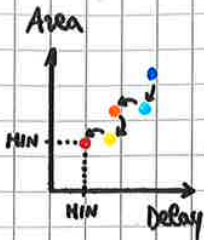
The cube is fully covered, therefore it is a tautology

Multi-level optimization and logic synthesis CH8

Summary: it is possible to implement a boolean function using a two-level representation (sum of products) or a multi-level representation.

The main differences between the two representations are the number of possible levels (2 vs ≥ 2) and the boolean gates used (AND, OR vs any restriction, as long as contained inside the tech library).

Motivations: why choose multi-level logic over two-level one? Because the implementation of circuits as multi-level logic networks allow **BETTER QUALITY** of final result.

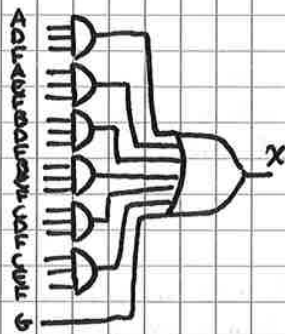


In a two-level optimization problem, area and delay are compatible objectives: you reduce area by minimizing the number of products, that allows for a smaller and moreover faster OR gate (lower # inputs \rightarrow lower delay). This is the **ONLY POSSIBLE SOLUTION!**

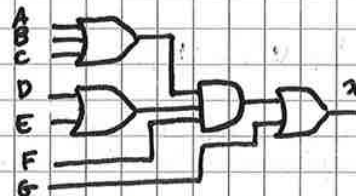
In a multi-level representation, you can explore the whole design space instead: hence you can implement a slow and small circuit up to a big but fast one. Moreover, usually the number of gates (as well as their # inputs) is reduced w.r.t a SOP implementation.

$$\text{Example: } x = ADF + AEF + BDF + BEF + CDF + CEF + G$$

$$x = (A+B+C)(D+E)F + G$$

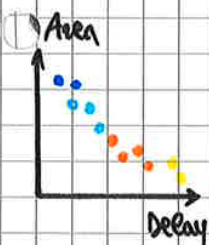


SOP representation



- 1) SMALLER AREA
- 2) REDUCE # INTERCONNECTIONS (\propto AREA)
- 3) MAYBE SLOWER

multi-level implementation



The ability of exploring the whole design space is the key for **Multiobjective optimization**: multiple constraints can be specified \rightarrow explore design space to find the best solution

The importance of multi-level synthesis grew in parallel with the growth of foundries for the semiconductor market: nowadays silicon vendors provides a pre-designed set of cells modeling basic logic (INV, NAND, NOR, ...) that is fully characterized (AREA, DELAY, POWER CONSUMPTION, ...) and optimized for the specific manufacturing process of that silicon vendor. The customer can not access the cell layout, and after synthesis and layout has a circuit ready to be manufactured.

Because transistor technology has become so complex (and expensive) the customer is not able to design its own cell; and because circuits are more and more complex he does not want to! This design approach is called **standard cell design**.

Estimators How to define metrics like area, delay, power consumption for a multilevel optimisation problem?

- Area: There is a direct proportionality between complexity and **# literals**. It is an easy and well accepted estimator.
- Delay: it depends on the set of vertices V_G (proportional to **# literals**) and set of edges E (that is estimated using MAX+SUM propagation explained later).
- Power: it is possible to define a probabilistic model based on switching activity at each node + capacitive loads (fanout).



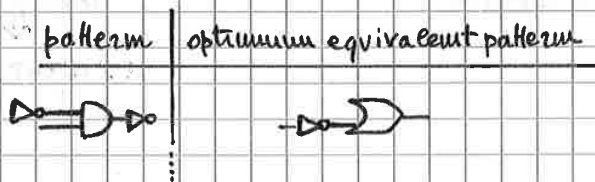
Methodology

The goal of multilevel synthesis is to obtain an equivalent representation of a given logic function that is optimal wrt a cost function.

With the topology-degree of freedom, even simplest problems becomes computationally hard. For this reason there are few exact methods (more of which is practice in real life) and some approximate optimization methods, classified in:

- Rule-based methods
- Heuristic algorithms

① Rule based methods rely on a rule database, composed by input patterns and an equivalent optimal implementation wrt a cost function.

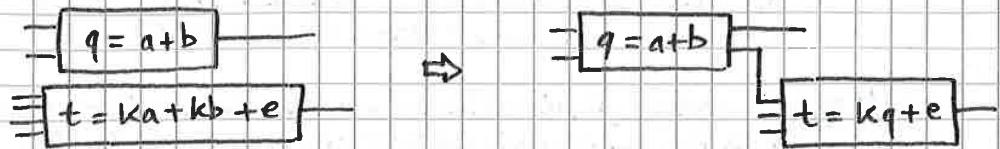


The method performs an algorithmic transformation, identifying patterns in the logic network and optimizing the design step by step.

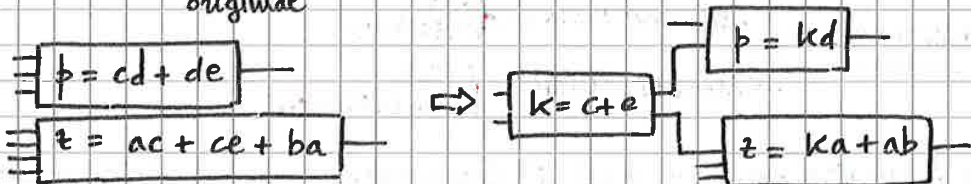
The two main drawbacks are:

- 1) BUILD AND **MAINTAIN** THE DATABASE, by selecting the most useful patterns for a particular design. Very time consuming and there is no evidence that it will affect positively the optimization.
- 2) **TECHNOLOGY CHANGE**: usually a new tech has different layouts to take into account the quantum effects: hence the previous may not apply the new tech \Rightarrow REBUILD DATABASE

Substitution: simplify a local function by using an additional input that was not previously in its support set. This transformation simplify local functions, but increase fanout



Extraction: this is one of the most powerful transformations; it finds a sub-expression of two (or more) expressions. It extract the subexpression, making a new node in the network and simplify the local functions.



Simplification: simplify a local function by means of an heuristic minimizer like Espresso. This may modify the fanin of the target node.



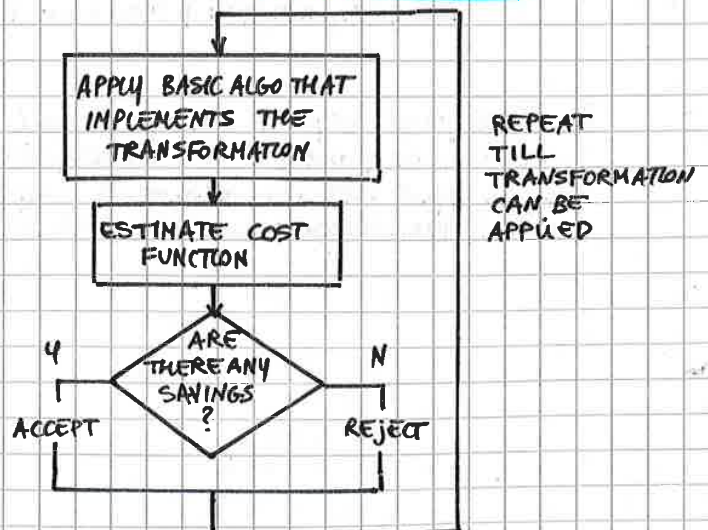
Algorithmic approach: implementation: we have seen what each transformation performs. But how a transformation is implemented?

Script composed of a sequence of transformations:

- TRANS 1 (DECOMPOSITION)
- TRANS 2 (...)
- TRANS 3
- ...
- TRANS N

COMPILE

HOW A TRANSFORMATION IS IMPLEMENTED?
AS AN OPERATOR



NB: some transformations require extra inputs (alongside the logic boolean network). Those inputs are the variables tuning the cost function.

Algebraic division

A function $f_{DIVISOR}$ is an algebraic divisor of another function $f_{DIVIDEND}$ if exist a function $f_{QUOTIENT}$ such that $f_{DIVIDEND} = f_{DIVISOR} \cdot f_{QUOTIENT} + f_{REMAINDER}$ if and only if:

1) $f_{DIVISOR} \cdot f_{QUOTIENT} \neq \{0\}$: this is trivial, because otherwise the result of the division is equal to $f_{REMAINDER} = f_{DIVIDEND}$. This condition must be verified also for the Boolean division.

2) $\text{sup}(f_{DIVISOR}) \cap \text{sup}(f_{QUOTIENT}) = \{0\}$ the intersection between the support sets must be void, otherwise some identity rules are needed (and are not allowed in Algebraic domain) - This property is not valid for Boolean division

Examples:

1) $f_{DIVIDEND} = ac + ad + bc + bd + e = (a+b)(c+d) + e$
 $f_{DIVISOR} = a+b$
 $f_{QUOTIENT} = c+d$
 $f_{REMAINDER} = e$

$\text{sup}(f_{QUOTIENT}) \cap \text{sup}(f_{DIVISOR}) = \{0\}$ OK!

2) $f_{DIVIDEND} = a + bc$
 $f_{DIVISOR} = a + b$
 $f_{QUOTIENT} = a + c$
 $f_{REMAINDER} = \{0\}$

$\text{sup}(f_{QUOTIENT}) \cap \text{sup}(f_{DIVISOR}) = \{a\}$ NOT OK

Let's check: $f_{DIVIDEND} = f_{DIVISOR} \cdot f_{QUOTIENT} + f_{REMAINDER}$
 $= (a+b)(a+c) = a + ac + ab + bc =$
 $= a(1 + b + c) + bc = a + bc$

↑
IT HAS NO MEANING IN ALGEBRAIC DOMAIN... (IDENTITY RULE)

↑
IN BOOLEAN DOMAIN

The Algebraic division is the basic operator used by 3 transformations:

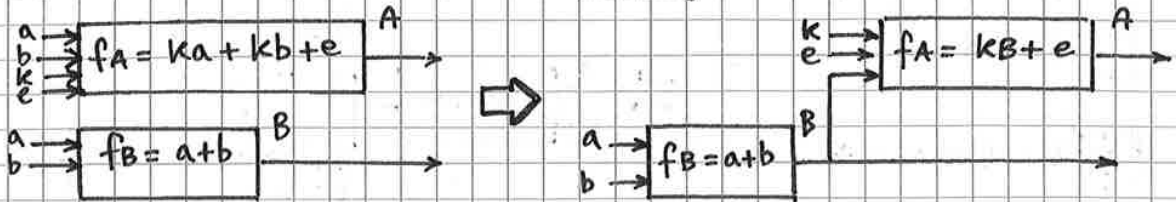
1) Decomposition: a mode (local function) is split into two smaller ones. The former mode is a divisor for the other mode

2) Substitution: a local function is simplified by using as additional input another mode of the network. The additional mode is a divisor for the original one

3) Extraction: a common sub-expression between two or more modes becomes a standalone mode. This new mode is a common divisor of the two original local functions

Algebraic substitution

The substitution transformation aims at simplify the local function of a given node by using an additional input previously not present in the support set. The additional input comes from an existing node of the network



The substitution, as can be easily understood, is a straightforward application of algebraic division. The algorithm can be summarized in the following steps:

- 1) Consider all possible couples of nodes
 - 2) Apply division (in any order) between them
 - 3) Use filters based on previous theorem to reduce computation
 - 4) If a quotient is not void, evaluate area and delay gain
 - 5) If they are higher than the user threshold, modify the network by increasing the fanout of one node (DIVISOR) and substituting fDIVIDEND with $f_{QUOTIENT} \cdot f_{DIVISOR} + f_{REMAINDER}$
- } complexity $\propto (\text{nodes})^2$
quadratic
- node label

A basic algorithm:

```

Substitution (G(V,E)) {
  foreach node i of network G {
    foreach node j of network G such that i ≠ j {
      A = set of cubes of fi
      B = set of cubes of fj
      if (A and B pass the filter test) {
        (Q,R) = Algebraic-division (A,B);
        if (Q ≠ ∅) {
          fQUOTIENT = sum of cubes of Q
          fREMAINDER = sum of cubes of R
          if (substitution is advantageous) {
            fi = j fQUOTIENT + fREMAINDER
          }
        }
      }
    }
  }
}
    
```


2) Recursive algorithm: it is based on the property that a "kernel of a kernel is still a kernel" (a kernel is of level "m" if has "m" recursive kernels)

two algorithms based on recursion are proposed:

R_KERNEL (with no pointers)

it does some redundant calculations

KERNEL (complete algo)

more optimized, exploit commutativity of mult.

Both algorithms uses a function CUBES(f, ci) which returns the cubes of f whose literals include those of cube ci

NB: kernel extraction is valid only if f is cube free \Rightarrow if f is not a kernel, we cannot extract kernels...

R_KERNEL(f) {

k = {} // kernel set initialized to void

foreach variable vi in sup(f) {

if (|CUBES(f, vi)| \geq 2) {

we are interested in cube free expressions, that implies multi-cube expressions \Rightarrow hence at least two cubes before performing division

Co-kernel \rightarrow C = maximal subcube among the cubes returned by CUBES(f, vi)

k = k \cup R_kernel(f/C)

RECURSION DIVISION!

}

k = k \cup f // the function itself is a kernel as well

Example:

f = ace + bce + de + g sup(f) = {a, b, c, d, e, g}

function: var	CUBES(f, var)	CO-K	K	
f: a	{ace} \geq 2 NO	/	/	
f: b	{bce} \geq 2 NO	/	/	
f: c	{ace, bce} \geq 2 YES	ce	R_kernel(f/ce)	recursion ev1
atb: a	{a} \geq 2 NO	/	/	
atb: b	{b} \geq 2 NO	/	/	k = {atb}
f: d	{de} \geq 2 NO	/	/	
f: e	{ace, bce, de} \geq 2 YES	e	R_kernel(f/e)	recursion ev1
actbctd: a	{ac} \geq 2 NO	/	/	
actbctd: b	{bc} \geq 2 NO	/	/	
actbctd: c	{ac, bc} \geq 2 YES	c	R_kernel(f/c)	recursion ev2
atb: a	{a} \geq 2 NO	/	/	
atb: b	{b} \geq 2 NO	/	/	k = k \cup {atb} = {atb}
actbctd: d	{d} \geq 2 NO	/	/	k = k \cup {atb, bctcd} = {atb, abtctd}
f: g	{g} \geq 2 NO	/	/	

Theorem (Brayton and McMullen)

Two expressions f_a and f_b have a common multiple-ube divisor f_d if and only if there exist at least 2 kernels, one from $K(f_a)$ and one from $K(f_b)$ for which $|K(f_a) \cap K(f_b)| \geq 2$. The common divisor to be extracted must be the one with the largest cardinality.

This theorem has two consequences:

- 1) It reinforces the idea that kernel set computation is necessary for the kernel intersection and hence the transformation extraction.
- 2) Introduces some filter conditions:
 - a) if an expression (local function) has no kernels, the corresponding mode can be dropped from consideration.
 - b) if kernel intersection is void, then the search for common sub-expressions can be dropped.

Finally, it is possible to define a basic algorithm for algebraic extraction, by combining kernel extraction, kernel set intersection and substitution:

Multiple-Extraction $(G(V, E), m, k) \{$

do {

foreach mode compute the kernel set of $e_{level} \leq k$
 compute all the kernel set intersections (Brayton theorem) and save them into a list L
 Rank the list L based on cardinality (descending order)

for $(i = 1 \text{ to } m) \{$

create a new mode j with local function $f_j = L(i)$ and add it to the network

} **ALGEBRAIC SUBSTITUTION** { Replace all modes whose local function has f_j as a divisor with \textcircled{j} QUOTIENT + REMAINDER **NO LABEL**

Update the kernel set of the modified modes by removing the kernel(s) in which the common divisor is contained

} **ALGEBRAIC SUBSTITUTION**

} while $(I \text{ is not empty})$

}

Final comments:

- 1) k represent the level for kernel extraction: basically, for each nested recursion, the level k increase. Hence this parameter set a bound on the complexity of kernel extraction.
- 2) m represent the amount of algebraic extractions performed before recomputing all kernel sets and kernel sets intersections. Hence it is a tradeoff between quality of result (small m) and execution speed (big m).

Example: $f = ace + bce + de + g$ $K(f) = \{ a+b, ac+bc+d, ace+bce+de+g \}$

decomposition is useless as the function itself.

1° approach: largest kernel

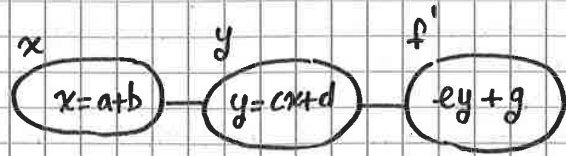
decomposition(f): $y = ac + bc + d$

$f' = ey + g$

recursion on y: $x = a + b$

$y' = cx + d$

$f'' = f' = ey + g$



2° approach: random kernel

decomposition(f): $y = a + b$

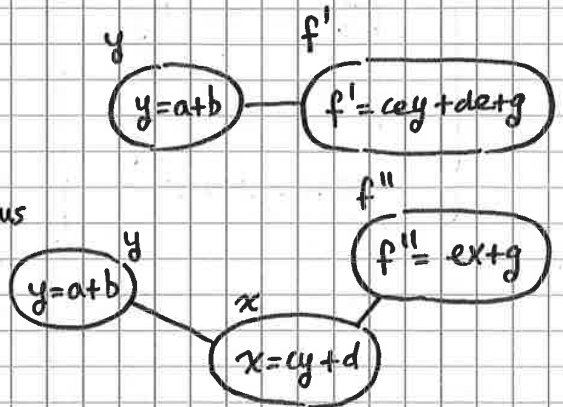
$f' = ce y + de + g$

recursion on y: no possible operations

decomposition(f'): $x = cy + d$

$f'' = ex + g$

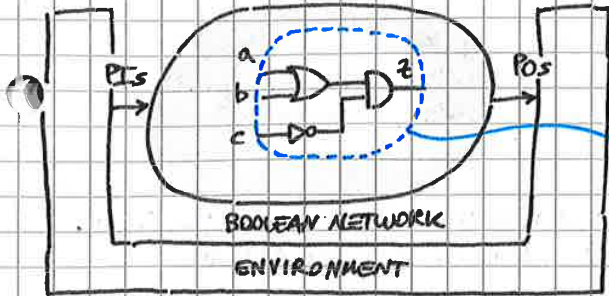
recursion on x: no possible operations



To achieve the same result, we need to run decomposition a larger number of times: more computationally expensive.

Summary on Algebraic methods: they abstract local functions as polynomials and perform division-based transformations. All transformations are reversible, hence it is easy to trade between area, delay and other metrics. They are fast and widely used methods, but they are weaker with respect to boolean methods.

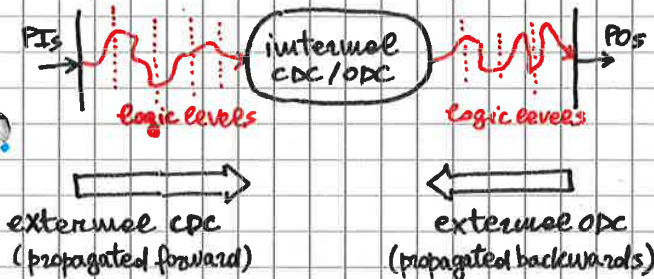
It is possible to generalize the concept of external don't cares:



it is possible to focus on a part of the boolean network, transforming the rest of the boolean network into the environment! It is then possible to define internal don't cares (ODC and CDC)

This approach allows to work with a much finer granularity, at the point that for each node of the network, an ODC and CDC set can be defined. Internal and external don't cares are different: who is the responsible?

The topology is the responsible; it may inject/remove don't care conditions by moving forward through the levels of the Boolean network:



For any generic node of the boolean network, the controllability of that node is defined by the external CDC propagated forward till that point. The observability of any node is given by the external ODC propagated backward till that point. Each logic gate on which ODC/CDC propagates may inject/remove new don't care conditions

Satisfiability don't care conditions:

The topology is responsible for the introduction/removal of internal don't care conditions. To find controllability don't cares generated by a local function (\equiv node), satisfiability don't care conditions must be used.

satisfiability don't care conditions:
$$SDC = \sum_{\text{all internal nodes}} x \oplus f_x$$

\downarrow OUTPUT AND NODE LABEL
 \uparrow LOCAL FUNCTION OF THE NODE

This simple condition represent all possible input/output combinations that are never produced \Rightarrow Hence don't care conditions. In fact, the XOR can be seen as the "opposite of a complement": each condition for which the two inputs (x and f_x) are different, makes the XOR true.

Example:



$$SDC_x = x \oplus (a' + b) = x'(a' + b) + x(a'b)$$

$$= x'(a' + b) + x(ab') = x'a' + x'b + xab'$$

DeMorgan = ab

Look closely: the first two conditions say that if one input is true, the output is false (never produced) the last condition says that the output is true if both inputs are false (never produced)

Example of iterative CDC computation

Compute CDC_{cut} knowing $CDC_{im} = \bar{x}_1 \bar{x}_4$

$$cut = \{PIs\} = \{x_1, x_2, x_3, x_4\}$$

$$CDC_{cut} = CDC_{im} = \bar{x}_1 \bar{x}_4$$

ITERATION 1

$$cut = \{x_1, x_2, x_3, x_4, a\}$$

$$SDC_a = a \oplus f_a = a \oplus (x_2 \oplus x_3)$$

$$CDC_{cut} = CDC_{cut} \cup SDC_a = \bar{x}_1 \bar{x}_4 + a \oplus (x_2 \oplus x_3)$$

$$D = \{x_2, x_3\}$$

$$P(CDC_{cut})|_{x_2} = \bar{x}_1 \bar{x}_4 + a \oplus (x_2 \oplus x_3) |_{x_2} \cdot \bar{x}_1 \bar{x}_4 + a \oplus (x_2 \oplus x_3) |_{\bar{x}_2} =$$

$$= (\bar{x}_1 \bar{x}_4 + a \oplus \bar{x}_3)(\bar{x}_1 \bar{x}_4 + a \oplus x_3)$$

$$P(CDC_{cut})|_{x_3} = [(\bar{x}_1 \bar{x}_4 + a \oplus \bar{x}_3)(\bar{x}_1 \bar{x}_4 + a \oplus x_3)] |_{x_3} \cdot [(\bar{x}_1 \bar{x}_4 + a \oplus \bar{x}_3)(\bar{x}_1 \bar{x}_4 + a \oplus x_3)] |_{\bar{x}_3}$$

$$= (\bar{x}_1 \bar{x}_4 + a)(\bar{x}_1 \bar{x}_4 + \bar{a})(\bar{x}_1 \bar{x}_4 + \bar{a})(\bar{x}_1 \bar{x}_4 + a)$$

can be dropped because resolution

$$= (\bar{x}_1 \bar{x}_4 + a)(\bar{x}_1 \bar{x}_4 + \bar{a}) = \bar{x}_1 \bar{x}_4 + a \bar{x}_1 \bar{x}_4 + \bar{a} \bar{x}_1 \bar{x}_4 = \bar{x}_1 \bar{x}_4$$

$$cut = cut - D = \{x_1, x_4, a\}$$

ITERATION 2

$$cut = \{x_1, x_4, a, b\}$$

$$SDC_b = b \oplus f_b = b \oplus (x_1 + a)$$

$$CDC_{cut} = CDC_{cut} \cup SDC_b = \bar{x}_1 \bar{x}_4 + b \oplus (x_1 + a)$$

$$D = \{x_1\}$$

$$P(CDC_{cut})|_{x_1} = [\bar{x}_1 \bar{x}_4 + b \oplus (x_1 + a)] |_{x_1} \cdot [\bar{x}_1 \bar{x}_4 + b \oplus (x_1 + a)] |_{\bar{x}_1}$$

$$= [0 + b \oplus 1] \cdot [\bar{x}_4 + b \oplus a] = \bar{b}(\bar{x}_4 + b \oplus a) =$$

$$= \bar{b}(\bar{x}_4 + \bar{b}a + \bar{a}b) = \bar{b}\bar{x}_4 + \bar{a}\bar{b}$$

$$cut = cut - D = \{x_4, a, b\}$$

ITERATION 3

$$cut = \{x_4, a, b, c\}$$

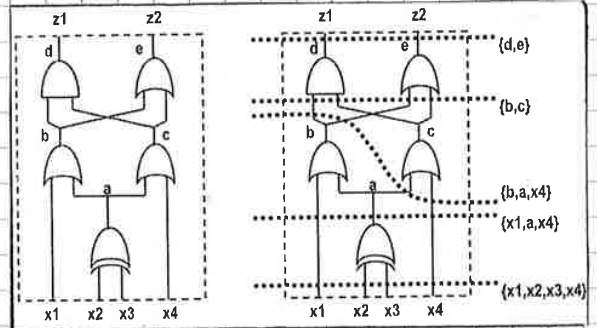
$$SDC_c = c \oplus f_c = c \oplus (a + x_4)$$

$$CDC_{cut} = CDC_{cut} \cup SDC_c = \bar{b}\bar{x}_4 + \bar{a}\bar{b} + c \oplus (a + x_4)$$

$$D = \{x_4, a\}$$

$$P(CDC_{cut})|_{x_4} = \bar{b}\bar{x}_4 + \bar{a}\bar{b} + c \oplus (a + x_4) |_{x_4} \cdot \bar{b}\bar{x}_4 + \bar{a}\bar{b} + c \oplus (a + x_4) |_{\bar{x}_4}$$

$$= [0 + \bar{a}\bar{b} + c \oplus 1] \cdot [\bar{b} + \bar{a}\bar{b} + c \oplus a] =$$



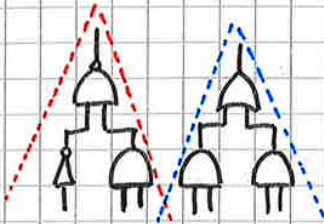
Internal ODC computation

The external ODC must be propagated backward among the logic levels of the Boolean network to compute the internal ODC at the desired mode. Each time the ODC goes through a logic level, don't care conditions can be added/removed. Hence two things are necessary:

1) A way to move backwards in the Boolean network: tree-network traversal algorithm. This time, there are two issues:

a) multiple outputs: if the Boolean network has multiple outputs, the problem became more complex, because observability depends on output: specific ODC conditions must be specified at each output, and one at a time they have to be propagated backwards.

This is the reason why during technology mapping (CHM) there is a phase called "partitioning" where the Boolean network is divided into "logic cones" (many-in, single-out network).



b) Non-tree network traversal: for networks with general topology, a node may have multiple fanout and multiple paths to any output. This situation, called fanout reconvergence is hard to handle both in logic optimization and testing problems. The problem lies in combining the observability of a variable along different paths. A naive, but wrong, assumption is that the ODC set in the intersection of the ODC sets related to the different fanout stems. Indeed, we must take into account the interplay of the ability of the paths to propagate to the output a change of polarity of the corresponding variable. More about reconvergence in testing and fault tolerance course.



2) A way to compute don't care conditions introduced by a mode of the Boolean network: the observability don't care set of a node is given by:

- the previous ODC conditions
- the complement of the boolean difference $\frac{\partial f_z}{\partial x}$, hence $f_z|_x \oplus f_z|\bar{x}$



In fact, the boolean difference represent if an output follows an input. Hence, the complement of the boolean difference represent all input combinations that do not modify the output \Rightarrow observability don't cares.

A possible algo could be:

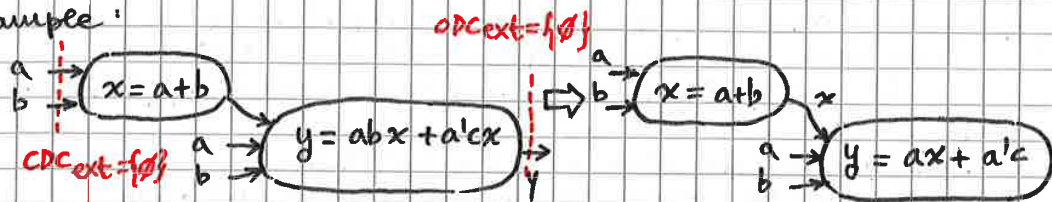
```

OBSERVABILITY (G(V,E), ODCout) {
  foreach mode x in reverse topological order {
    ODC = ODC +  $\frac{\partial f_{current mode}}{\partial x}$ 
  }
}
    
```


Boolean simplification

At the beginning of this chapter, it was said that all 5 algebraic transformations can be collapsed into one boolean transformation: the simplification. The goal is to manipulate a local function to obtain one with less literals by using the internal don't care conditions. The difference with respect to heuristic minimizers (like Espresso) is that the boolean simplification uses don't care conditions coming from the nodes around (topology), while Espresso can only rely on don't care conditions of the single node.

Example:



Suppose external CDC and external ODC are empty. The CDC and ODC conditions can be propagated through the network:

$$CDC_x = x \oplus fx = x \oplus (a + b) = ab'x + a'x' + bx'$$

$$ODC_x = \frac{\partial fy}{\partial x} = (abx + a'cx) \Big|_x \oplus (abx + a'cx) \Big|_{\bar{x}} = 1 \oplus 0 = 0$$

$$DC_x = CDC_x + ODC_x = ab'x + a'x' + bx'$$

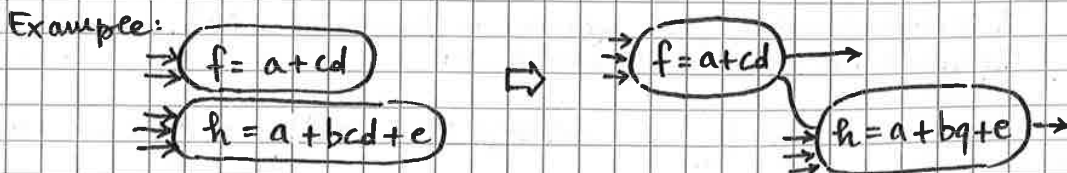
By exploiting don't care conditions, the node y can be simplified:

$$y = ax + a'c$$

because ax never produced because $a'x'$ is never produced

Boolean substitution

The substitution is empowered by the boolean simplification: given two nodes, the support set of one must be modified to include the other node; then, by means of don't care condition, the node can be simplified, otherwise, if no optimization is possible, the famous connection is removed.



$$SDC_f = f \oplus f_f = f \oplus (a + cd) = \bar{f}(a + cd) + f(\overline{a + cd}) = a\bar{f} + \bar{f}cd + f\bar{a}\bar{c}\bar{d}$$

is $cd = f$? this relation is true when

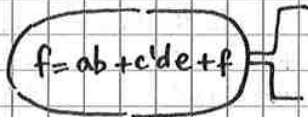
$$\begin{cases} cd_f \\ \bar{c}\bar{d}_f \end{cases} \in DC_{set} \text{ in the DC set there is } q\bar{c}\bar{d}\bar{a} : \text{ when } a \text{ is true, } f \text{ is true, so it doesn't matter}$$

Timing issues in multi level logic optimisation CH10

Delay is a crucial parameter of an electronic circuit: during optimisation it is possible to either select it as a constraint and search for minimum area, or vice versa set an Area constraint and search for minimum delay. Delay estimation is important also during verification, to verify the I/O delay constraints as well as cycle time constraints. Hence it is crucial to define a delay estimator / timing model.

Delay modeling: it is possible to distinguish:

- 1) Logic / stage (gate) delay modeling: given a generic Boolean network, each node represents a logic function. Depending on the local function complexity, the delay may change. A good estimator takes into account the # literals as well as the node fanout:

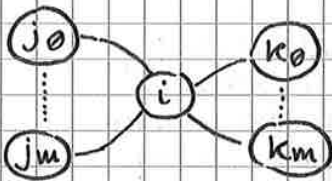


literals = 6
fanout = 2
Delay = $\alpha \cdot 6 + \beta \cdot 2$

Delay = $\alpha \cdot \# \text{ literals} + \beta \cdot \text{fanout}$

where α and β are weighting coefficients. A cheaper estimator can take into account only the number of literals. Each node of the network is annotated with this delay estimation. In case of mapped Boolean networks, delay can be found in tech library.

- 2) Network delay modeling: for each node it is possible to define:



- a) Propagation delay: a positive number, possibly accounting for load or function size (computed with gate delay modeling). Notation: d_i
- b) Arrival time: it is defined at the input of the considered network, and can be propagated forward through the logic levels:

$$t_i = \max(t_j | e_j \rightarrow i \in E) + d_i$$
NODE DELAY

the arrival time is the time after which the output of the node is ready

- c) Required time: it is specified at primary outputs, and can be propagated backwards through the logic levels:

$$\bar{t}_i = \min(\bar{t}_k - d_k | e_i \rightarrow k \in E)$$

- d) Slack: difference between required time and arrival time:

$$s_i = \bar{t}_i - t_i$$

If positive, there is some margin
If negative, there is a violation

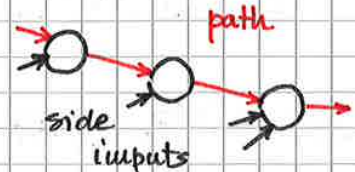
Sensitizable path

How to check if a topological critical path is a true critical path?

If a path is ~~sensitizable~~ sensitizable, an event can propagate from its tail to its head. Two definitions of sensitization exist:

- 1) **Dynamic sensitization:** for each node of the path, it must be verified if the output follows the input transition. This condition can be easily modeled with **boolean derivative**. The boolean derivatives must be true **at the time at which the event propagates**. The boolean derivatives will be function of side inputs

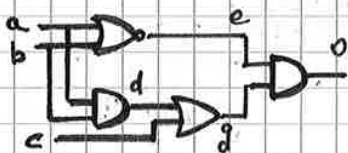
$$\frac{\partial f v_i}{\partial v_{i-1}} = 1 @ t_{v_{i-1}} \quad \forall v_i \in T.C.P.$$



- 2) **Static sensitization:** like for dynamic sensitization, for each node of the path it must be verified that the output follows the input transition. This condition is modeled with boolean derivative. The boolean derivatives will be function of side inputs. This time, the timing information is neglected: this choice makes the model simpler but weaker. In fact it may lead to **underestimated delays** due to glitches or fanout reconvergence.

$$\frac{\partial f v_i}{\partial v_{i-1}} = 1 \quad \forall v_i \in T.C.P.$$

Example:



Suppose all gates have unit delay. Then the topological critical paths are:
 $\{a, d, g, o\}$, $\{b, d, g, o\}$
 Due to circuit symmetry, we can analyze only one.

STATIC

$$\frac{\partial f v_o}{\partial g} = e \cdot g|_g \oplus e \cdot \bar{g}|_{\bar{g}} = e \oplus \emptyset = e$$

$$\frac{\partial f v_g}{\partial d} = d \cdot c|_d \oplus d \cdot \bar{c}|_{\bar{d}} = 1 \oplus c = \bar{c}$$

$$\frac{\partial f v_d}{\partial a} = a \cdot b|_a \oplus a \cdot \bar{b}|_{\bar{a}} = b \oplus \emptyset = b$$

The three conditions must be valid (=1) together: therefore $AND(\text{conditions}) = 1$

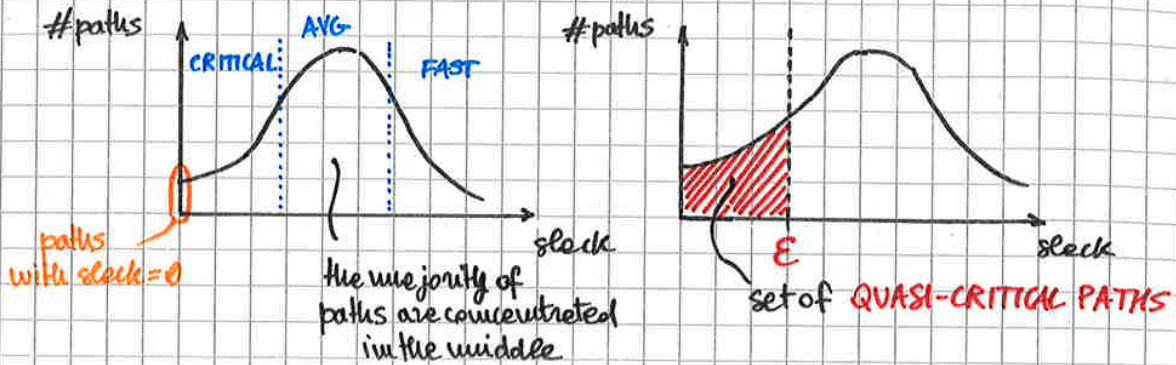
$$e \cdot \bar{c} \cdot b = 1 \rightarrow (\bar{a} + b) \bar{c} \cdot b = \bar{a} \bar{c} b = 0$$

This path is not statically sensitizable.

Algorithm for delay optimization/minimization

Given as inputs the arrival time of PIs and required time of POs, the algorithm must alternate between critical path computation and logic transformations on nodes belonging to critical paths.

A possible algorithm must iterate on critical paths: but how to discriminate between "normal" and "critical" paths? The user must specify a slack threshold ϵ :



The algorithm must operate on multiple paths at the same time to:

- 1) Improve quality of result
- 2) Reduce number of transformations applied
- 3) Avoid "ping-pong" effect
- 4) Improve convergence of the algorithm

A possible algorithm is the following:

```

Reduce_Delay ( G(V,E), ti of PIs,  $\bar{t}_i$  of POs,  $\epsilon$  ) {
  repeat {
    propagate the arrival time forward
    propagate require time backward
    compute the slack of each node
    Identify topological critical paths
    Filter TCPs using sensitization conditions

    Identify all sensitizable paths with slack lower than  $\epsilon$ 

    Select a subset of the total nodes on which apply transformations
  } until (no possible savings)
}
    
```

static timing analysis

Algorithmic based approach

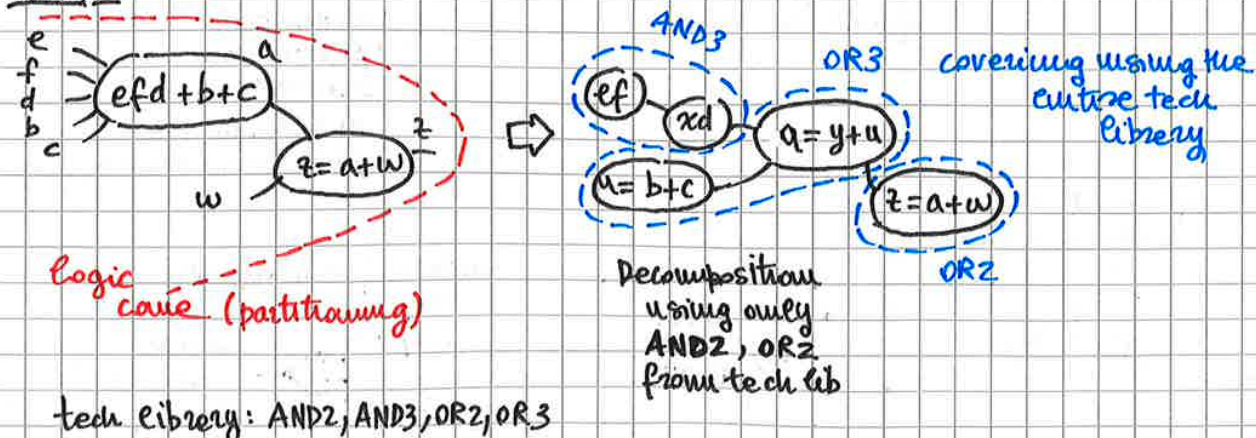
The algorithmic-based approach can be subdivided in 3 substages:

- 1) PARTITIONING: this stage aims at dividing the unbounded generic boolean network (multi input, multi output) into a set of logic cones (multi input, single output).
- 2) DECOMPOSITION: each logic cone is composed by generic (and maybe complex) local functions. This stage decomposes the logic cones (using decomposition transformation) using a subset of Boolean operators available in the library (in our course we use NAND-INV). This step is crucial because:
 - a) guarantees the existence of a solution (trivial binding)
 - b) simplifies the algorithmic solution of the covering problem
- 3) COVERING: this stage aims at optimizing the trivial binding solution found at the decomposition stage. To achieve this goal, two main operations must be performed:
 - a) identify covers: partition the logic cone into a set of smaller Boolean networks; each one of them can be replaced with a specific gate inside the tech library
 - b) matching: a smaller Boolean network and a gate (cell) in tech library do match if they show the same terminal behavior.

The covering problem can be solved by means of:

- a) Exact methods based on Boolean matching
- b) Heuristic methods (dynamic programming) based on Boolean tautology or structural approach (tree based).

Example:



A possible algorithm for covering (tree based):

for each node in SUBJECTTREE (from input to output) {

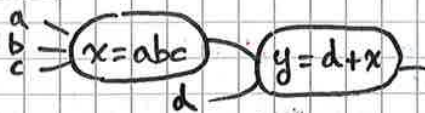
- 1) Check for possible matching. There are two possible options:
 - a) the locally rooted subtree (till the inputs) does match with a pattern tree (cell)

⇒ COST OF CELL
 - b) a subtree of the locally rooted subtree does match with a pattern tree (cell)

⇒ COST OF CELL + COST OF LEAVES
- 2) take the one with minimum cost

}

Example:



The given tech library is composed of

INV	\bar{a}	cost = 1
NAND2	\overline{ab}	cost = 2
AND2	ab	cost = 4
NOR2	$\overline{a+b}$	cost = 3
OR2	$a+b$	cost = 5
AOI21	$\overline{ab+c}$	cost = 7
AOI22	$\overline{ab+cd}$	cost = 8
NAND3	\overline{abc}	cost = 4

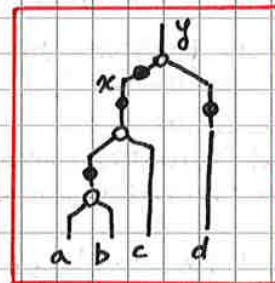
1° stage: partitioning. The boolean network is already a logic cone with only one output

2° stage Decomposition: bind the network to NAND2 + INV cells:

$$x = abc = \overline{\overline{abc}} = \overline{\overline{ab} \cdot c}$$

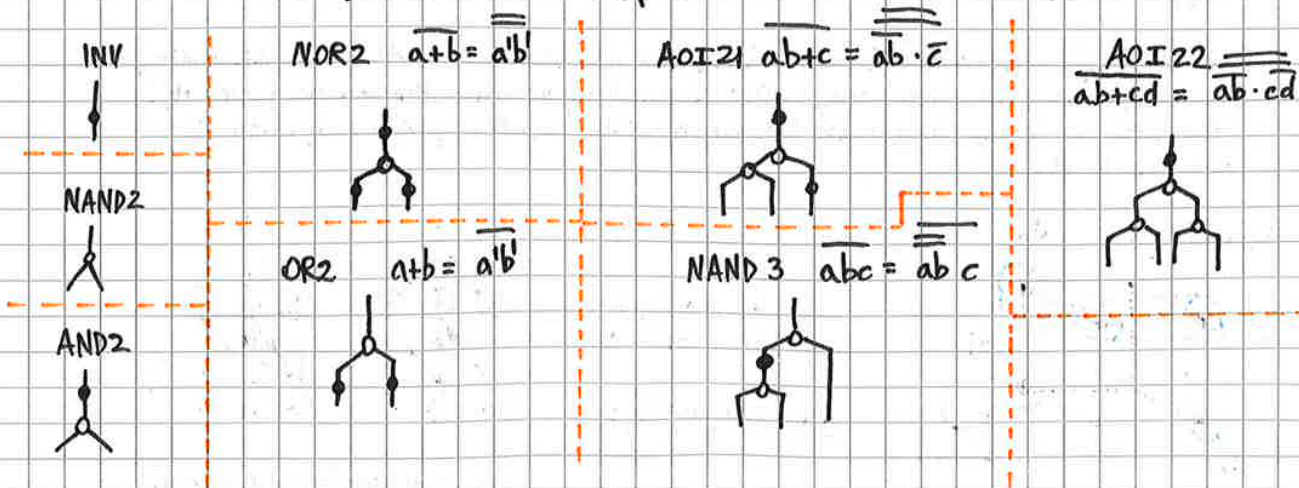
$$y = d+x = \overline{\overline{d+x}} = \overline{\overline{d} \cdot \overline{x}}$$

Let's identify NAND2 as and INV as



SUBJECT TREE

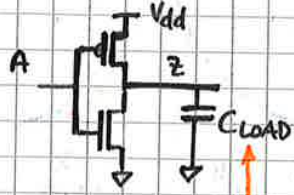
3° stage: Covering. Build the network using the entire tech library. Therefore, we must build first the library as tree patterns using only NAND2 and INV. In this way, possible matching on subject tree can happen.



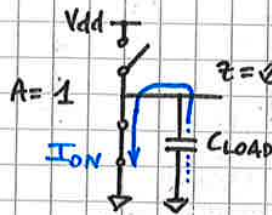
Power consumption modeling and estimation CH 12

Earlier in the course an area estimator (# literals, or if you have a tech library, sum of physical area occupied by each gate of logic circuit) and a delay estimator (based on static timing analysis engine, also implemented in commercial tools) were introduced. What is missing is a power estimation(s).

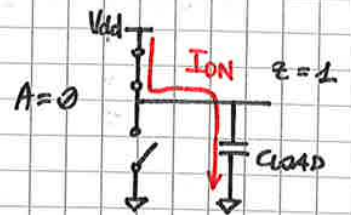
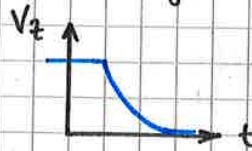
For a CMOS circuit, the power consumptions can be divided into 3 main contributions: dynamic, static and short circuit power. Let's use an example: a CMOS inverter:



LUMPED CAPACITANCE
taking into account:
- cap. of fanout gates
- cap. of wires
- etc.



capacitance is discharged



capacitance is charged



$$\text{Energy} = C_{LOAD} \cdot V_{DD}^2$$

$$\text{Energy STORED} = \frac{1}{2} \text{Energy} = \frac{1}{2} C_{LOAD} V_{DD}^2$$

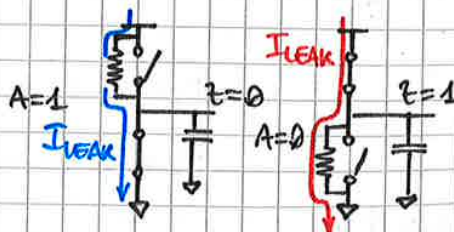
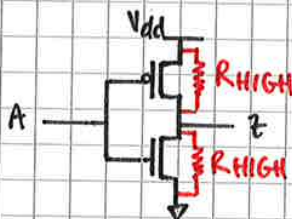
$$\text{Energy WASTED} = \frac{1}{2} C_{LOAD} V_{DD}^2 \quad \text{energy consumed as HEAT}$$

$$P = \frac{\Delta E}{\Delta T} = \frac{1}{2} C_{LOAD} V_{DD}^2 \frac{1}{T_{CLK}} = \frac{1}{2} C_{LOAD} V_{DD}^2 f_{CLK}$$

$$P_{dynamic} = \frac{1}{2} C_{LOAD} V_{DD}^2 f_{CLK} \cdot E_{sw}$$

where E_{sw} is the switching activity: it represents the probability that a given gate output makes a transition

Due to non idealities in the CMOS transistors, they cannot be modeled as ideal switches. Even if they are open, there is still a leakage current flowing, that is modeled by means of an High value resistor R_{HIGH} :



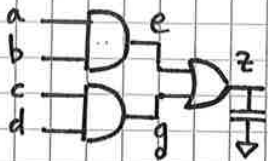
$$P_{static} = I_{leakage} \cdot V_{DD}$$

Example : state client



$$\left. \begin{aligned} P_z: 0 \rightarrow 1 &= P_{\bar{z}} P_z \\ P_z: 1 \rightarrow 0 &= P_z P_{\bar{z}} \end{aligned} \right\} E_{sw} = P_z: 0 \rightarrow 1 + P_z: 1 \rightarrow 0 = 2 P_z P_{\bar{z}} = 2 P_z (1 - P_z)$$

Example : switching activity propagation over a given topology



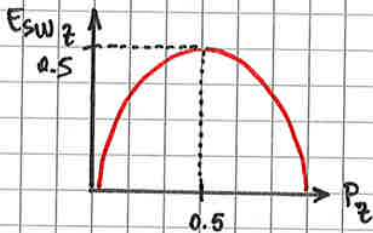
Suppose
 $P_a = P_b = P_c = P_d = \frac{1}{2}$

$$P_e = P_a \cdot P_b = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} \rightarrow P_{\bar{e}} = 1 - P_e = \frac{3}{4} \quad E_{sw_e} = 2 P_e (1 - P_e) = 2 \cdot \frac{1}{4} \cdot \frac{3}{4} = \frac{3}{8}$$

$$P_g = P_c \cdot P_d = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} \rightarrow P_{\bar{g}} = 1 - P_g = \frac{3}{4} \quad E_{sw_g} = 2 P_g (1 - P_g) = 2 \cdot \frac{1}{4} \cdot \frac{3}{4} = \frac{3}{8}$$

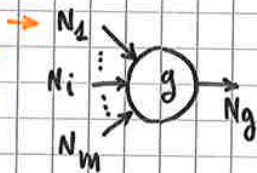
$$P_z = 1 - P_e P_g = 1 - \frac{3}{4} \cdot \frac{3}{4} = \frac{7}{16} \rightarrow P_{\bar{z}} = 1 - P_z = \frac{9}{16} \quad E_{sw_z} = 2 P_z (1 - P_z) = \frac{7}{16} \cdot \frac{9}{16} = \frac{63}{128}$$

If the static probability is not propagated through the circuit, the switching activity for each node cannot be computed. Therefore Synopsys best to assume the worst case scenario: $P_z = \frac{1}{2} \rightarrow E_{sw} = \frac{1}{2}$



Moreover, the switching activity alone is not enough for considering the dynamic behavior of the circuit: the **toggle number** must be computed:

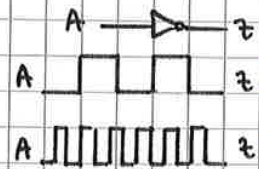
NUMBER OF TOGGLES OF INPUT



NUMBER OF TOGGLES

$$N_g = \sum_{\text{inputs}} N_i P \left(\frac{\partial f_g}{\partial i} \right)$$

BOOLEAN DIFFERENCE



same E_{sw} , different power consumption

Example:

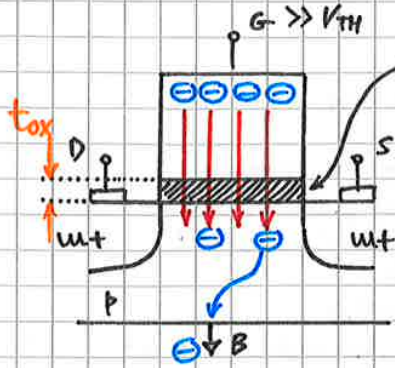


$$N_z = N_A \cdot P \left(\frac{\partial f_z}{\partial A} \right) + N_B \cdot P \left(\frac{\partial f_z}{\partial B} \right)$$

$$\frac{\partial f_z}{\partial A} = ab|_A \oplus ab|_{\bar{A}} = b \oplus 0 = b$$

$$\frac{\partial f_z}{\partial B} = ab|_B \oplus ab|_{\bar{B}} = a \oplus 0 = a$$

$$N_z = N_A \cdot P(b) + N_B \cdot P(a)$$

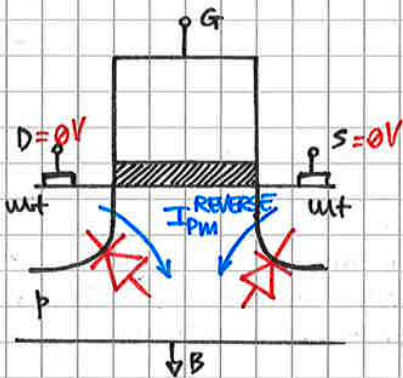
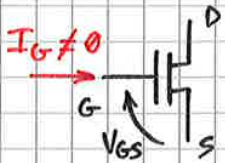


SiO₂: silicon dioxide (insulator)

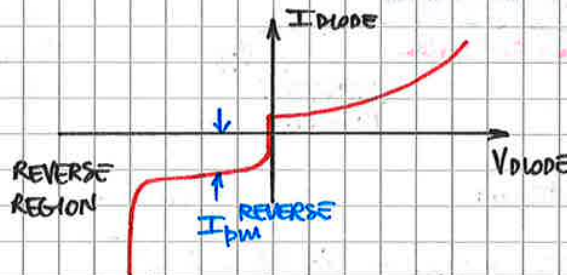
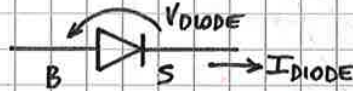
The gate metal and the p-doped semiconductor have in between a thin layer of insulator: this is a parallel plate capacitor!

Therefore, when V_G is positive, a strong electric field is generated: due to quantum effects (tunneling), some electrons overcome the insulator voltage barrier, reaching the p-doped semiconductor, and flowing into the bulk and so into ground - To summarize, due to quantum effects, there is a gate current non negligible =

$$I_{GATE} \propto WL \left(\frac{V_{GS}}{tox} \right)^2$$



The n-doped and p-doped regions form a p-n junction \equiv DIODE!
If the bulk voltage becomes sufficiently large, a current will flow

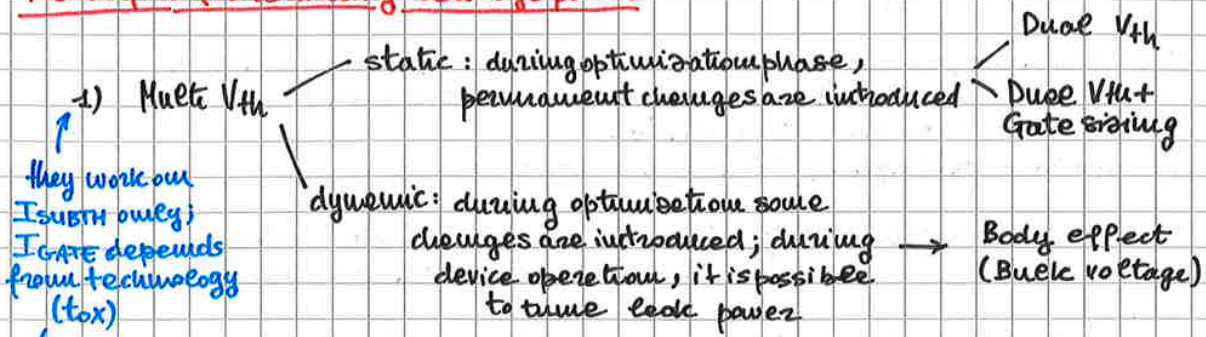


To summarize

$$I_{LEAKAGE} = I_{SUBTH} + I_{GATE} + I_{PN}^{REVERSE}$$

$$P_{STATIC} = I_{LEAKAGE} \cdot V_{dd}$$

Techniques for reducing leakage power



2) Input based: a certain circuit state consume less leakage power than another.

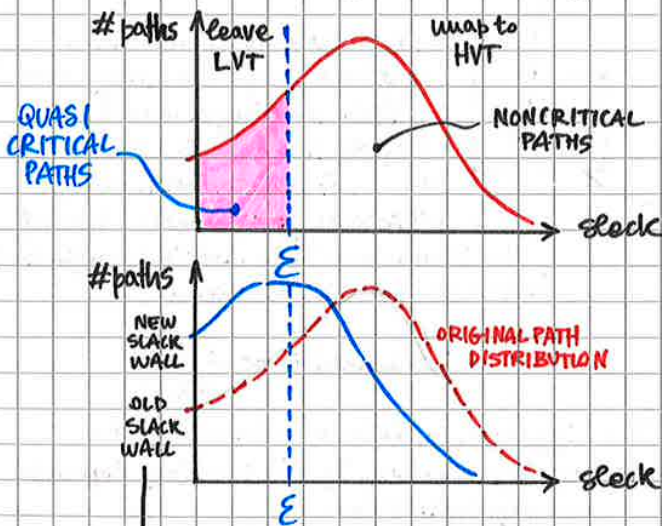
Let's talk about Multi V_{th} (Dual) ^{static} optimization: usually a standard cell library is "duplicated" for different threshold voltages.

The library used during laboratories experiences contains:

HVT (high V_{th})	leakage ↓	Propagation delay ↑
SVT (standard V_{th})	—	—
LVT (low V_{th})	leakage ↑	Propagation delay ↓

Therefore, low threshold voltage cells are faster but leaky, while high threshold voltage cells are slower but less leaky with respect to the standard V_{th} cells.

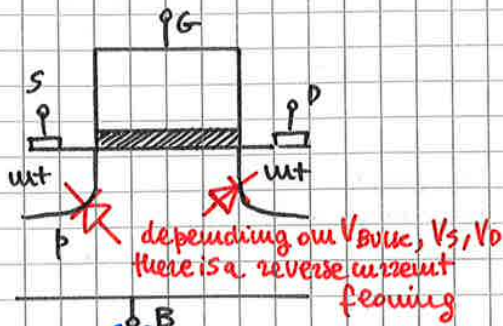
The multi V_{th} approach allows leakage power minimization WITHOUT reducing the operating frequency:



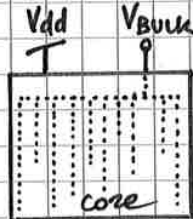
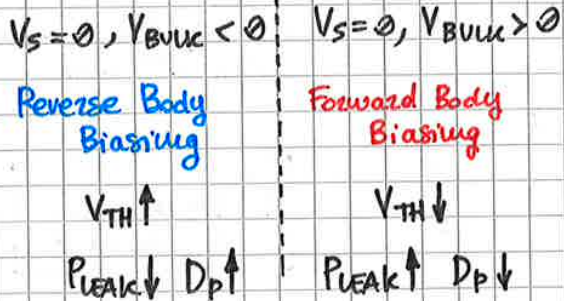
- 1) Synthesize circuit using only LVT cells to obtain best performances.
- 2) Identify a threshold ϵ for slack, to distinguish between quasi-critical paths and non-critical paths.
- 3) Swap cells of non-critical paths with HVT cells, to keep same overall timing performances but less power consumption (leakage)

Slack wall: the number of paths with slack = 0 increases. Due to process variations, the probability that a slower cell is in a timing critical path increases, and therefore the probability that the manufactured product works are decreased (\equiv yield ↓)

Multi V_{th} dynamic approach: Body effect



V_{BULK} CAN BE TUNED DYNAMICALLY AT RUN TIME



During IDLE, the bulk voltage is driven to a very low value
 ⇒ **STAND BY LEAKAGE**

The core is slower, but because it is not used, it does not matter!

Another application for this technique:

Suppose that after manufacturing, process variations made your device too slow: instead of decreasing the clock frequency or binning the device, the V_{BULK} can be increased such that D_p is compliant with specs. This can effectively increase yield, but V_{BULK} can be increased by a small amount or otherwise too much leakage power.

Input based approach for reducing leakage

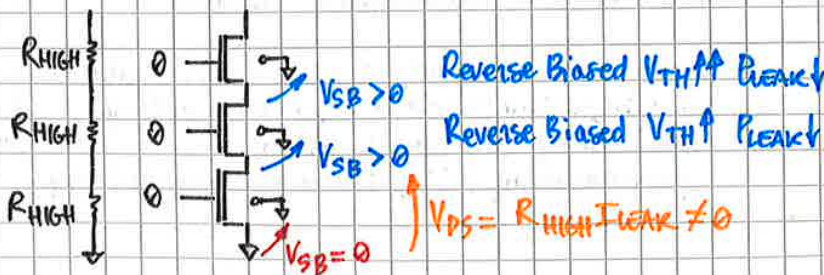
This technique is very easy to implement, but is not very effective as the previous ones. We already know that depending on the input state, a gate can consume a different amount of leakage power (see lib files).

For example, it is possible to seek input combinations (based on P_{LEAK}) of a CMOS NAND gate:

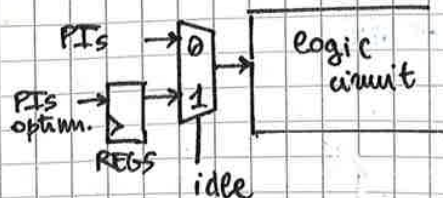
A	B	Z	Leak
0	0	1	1 (MINIMUM LEAKAGE)
0	1	1	3
1	0	1	2
1	1	0	4 (MAX LEAKAGE)

This phenomenon is related to body biasing - In fact a chain of NMOS is more leakage efficient than a parallel of PMOS (in this specific case of the NAND GATE)

eq. model

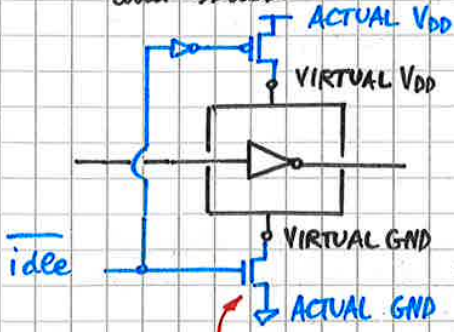


Therefore, there is an optimal input pattern such that leakage is minimized:



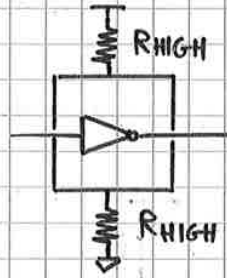
Power gating

This approach is in the middle between static and dynamic.
It allows 100% leakage savings, but has some drawbacks and issues.



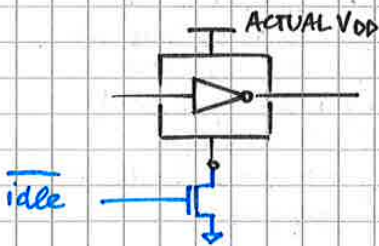
SLEEP TRANSISTOR: HVT cell
with almost no leakage
(very high V_{TH} ...)

in reality



power gating can be applied only to functional cores. It is possible to tune the sleeping strength.

The same result can be obtained with one transistor (NMOS) only!
In this way 50% area is saved, but some issues may arise... I'll discuss it later on

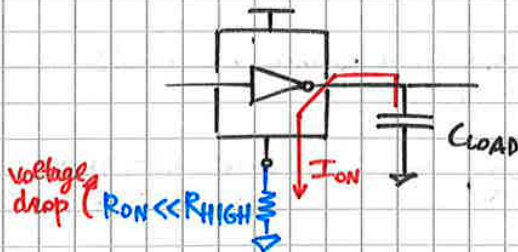


One of the most noticeable drawbacks is that the cell becomes slower!
The sleep transistor may become a bottleneck for current I_{DS} :
it must be sized properly:

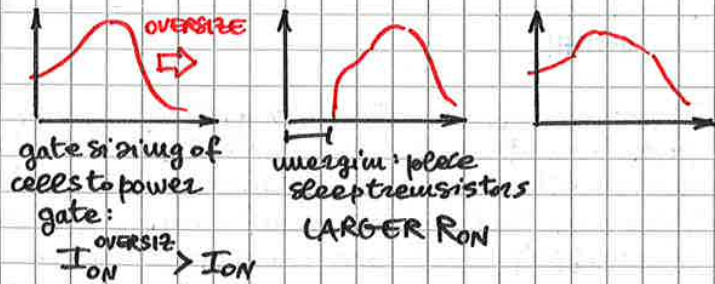
$$\Delta D_p^{MAX} \leq 3\% D_p^{NOMINAL}$$

$$R_{ON}^{MAX} = \frac{V_X^{MAX}}{I_{ON}}$$

To overcome the propagation delay issue, you must oversize the cells and use the slack margin for the sleep transistors



voltage drop ($R_{ON} \ll R_{HIGH}$)



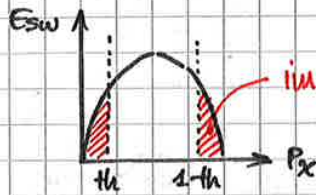
How to find out the maximum I_{DS} current to correctly size the sleep transistor? Very complex to calculate: try all possible permutations of all possible inputs patterns, run a full simulation and annotate biggest current. Too much consuming, use Heuristics!

Accounting for signal probabilities

Given a gate and probability of its inputs, the switching activity can be defined as:

$$E_{sw} = 2 P_x (1 - P_x)$$

where P_x is the probability that output x of considered gate is true. There is an interesting correlation between P_x and E_{sw} :



in the highlighted regions the switching activity is low. Therefore during optimization, a solution that provides a smaller E_{sw} can be preferred to another one, by modifying the cost function/estimators:

$$\text{area} = \# \text{ literals} \cdot P_x (1 - P_x)$$

Or an heuristic can be exploited:

- 1) set probability threshold th
- 2) for each node compute don't care set and output signal probability
 - a) if probability is a desired one (below th or above $1 - th$) minimize the function using don't cares and traditional cost function
 - b) if probability is not a desired value minimize the function using don't care conditions and modified cost function that takes into account for signal probability

State Assignment

In an FSM, to minimize power consumed by state registers, the binary code of each state must be assigned taking into account transition probabilities: the larger the prob, the smaller the hamming distance between the two state codes:

