



Compilers

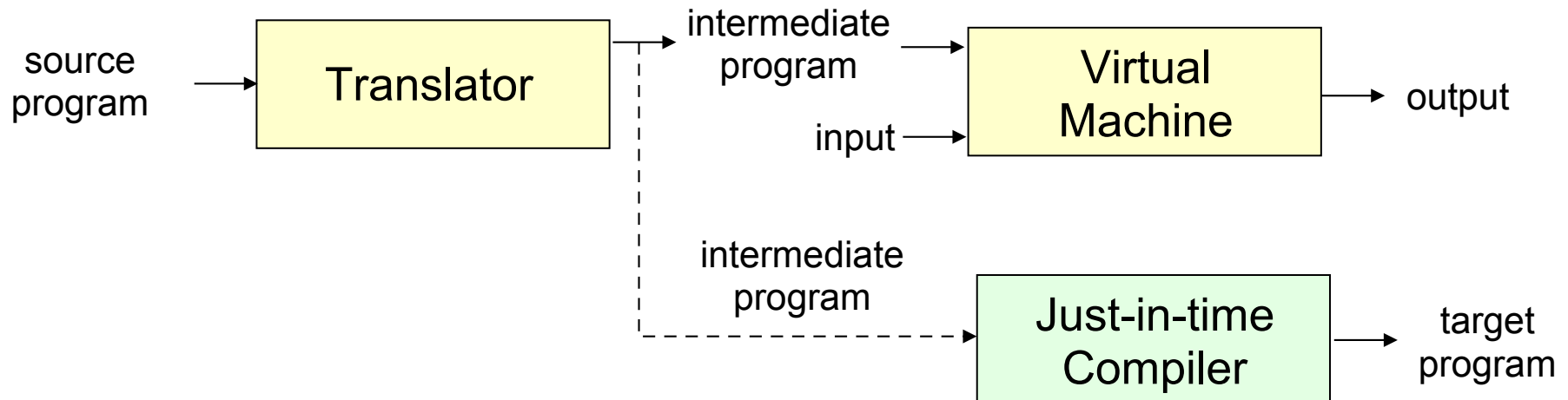
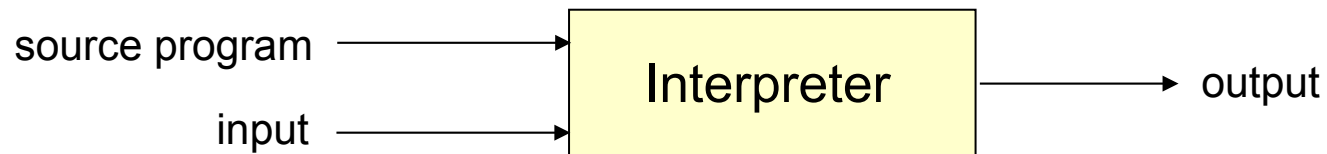
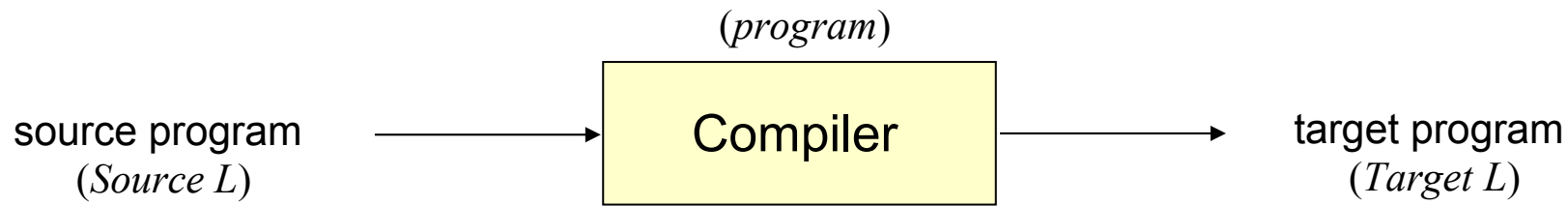
1	Introduction
2	Lexical analysis
3	Syntax analysis
4	Semantic analysis
5	Code generation
6	Runtime environments

- Aho, Lam, Sethi, Ullman “*Compilers. Principles, Techniques, and Tools*”, Addison-Wesley, 2006
- Levine, Mason, Brown “*Lex & Yacc*”, O’Reilly & Associates.
- Slides : <http://www.ing.unibs.it/lamperti>

History Background

- Link $\left\langle \begin{array}{l} \text{Programming languages} \\ \text{Implementation technology} \end{array} \right\rangle \Rightarrow$ abstraction process possible with compilers
- Late 1940s: first computers based on the architecture of von Neumann
→ need for writing programs = [*statements*]
- Machine code: `C7 06 0000 0002` = hexadecimal statement of Intel 8x86 for inserting number 2 at address 0000
- Assembly language: `MOV X, 2` \Rightarrow translated into machine code by an assembler  $\left\langle \begin{array}{l} \text{difficult to understand} \\ \text{not portable} \end{array} \right\rangle$
- “Ideally”: specification of program operations in concise form $\left\langle \begin{array}{l} \text{mathematics (expressions)} \\ \text{natural language (control)} \end{array} \right\rangle$
`X = 2` 
- Historic objections $\left\langle \begin{array}{l} \text{impossible} \\ \text{if possible} \rightarrow \text{inefficient target code} \end{array} \right\rangle \Rightarrow$ contradicted by **FORTRAN**

Language Processors



- Also: source-to-source translators (e.g.: C++ → C)

Language Processors (ii)

- Apparent complexity factors in designing general techniques for compilers:
 1. Thousands of source L
 2. Thousands of target L $\left\{ \begin{array}{l} \text{programming L (source-to-source)} \\ \text{machine L} \end{array} \right.$
 3. Various compiler typologies (single-pass, multiple passes, different grouping of phases ...)
- At a certain abstraction level → complexity reduced to a few essential procedures
→ make use of the same techniques
- Techniques for compilers → reusable in other contexts of computer science
- Realms of interest: $\left\{ \begin{array}{l} \text{Artificial-language theory} \\ \text{Programming languages} \\ \text{Web} \\ \text{Hardware architectures} \\ \text{Algorithms} \\ \text{Software engineering} \end{array} \right.$

Implementation of High-Level Programming Languages

- High-level programming languages: easier to program but less efficient
- Low-level programming languages: more efficient but harder to write, less portable, more prone to errors, harder to maintain
- Optimizing compilers → techniques to improve efficiency of generated code
- Example: **register** keyword in C: considered necessary (in mid 1970s) to control which variables reside in registers
 - No longer necessary with effective register-allocation techniques
 - Hardwiring register allocation may hurt performances!

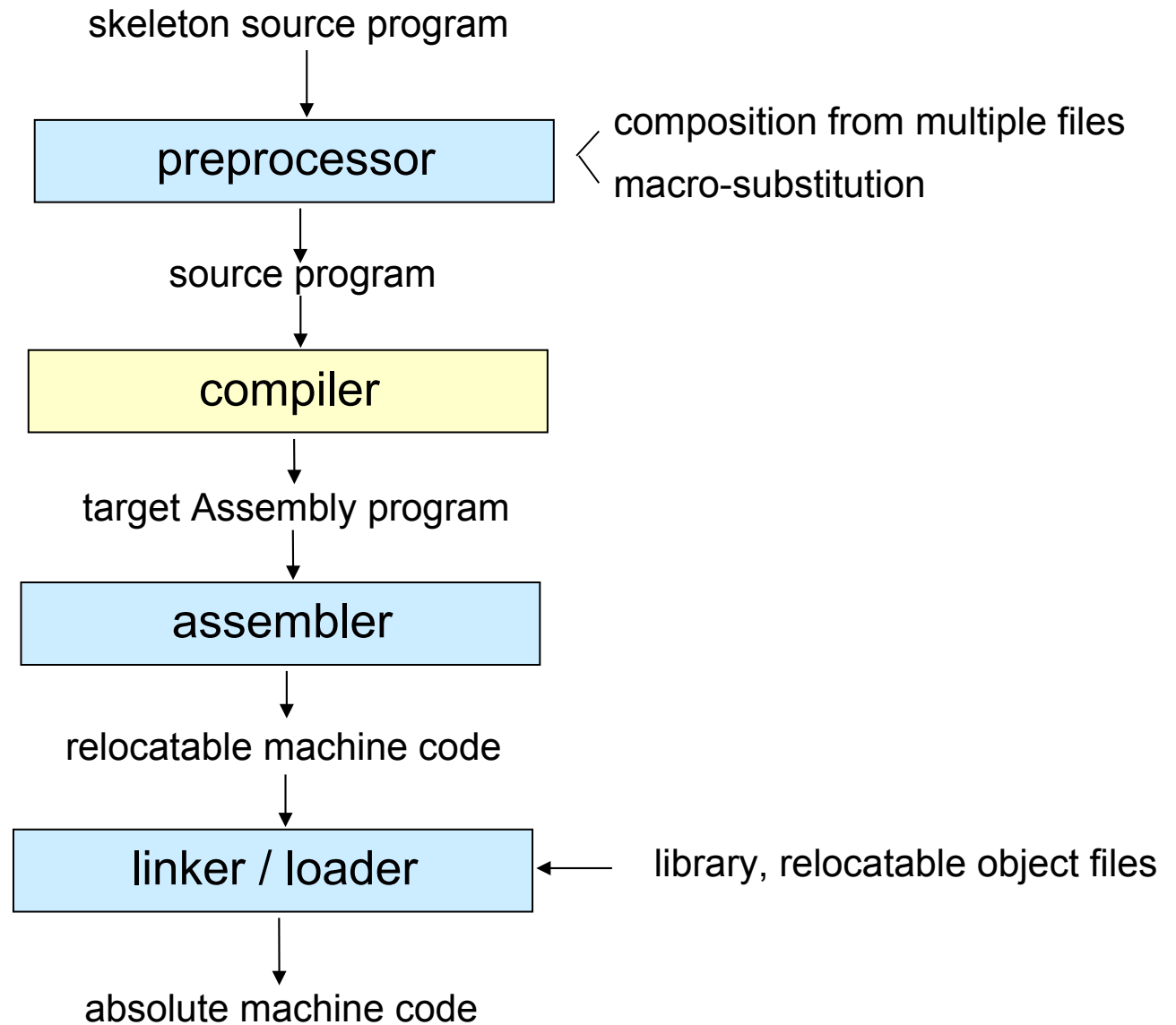
Optimizations for Computer Architectures

- Rapid evolution of computer architecture → demand for new compiler technology
- High-performance systems take advantage of $\begin{cases} \text{parallelism} \\ \text{memory hierarchies} \end{cases}$
- **Parallelism:**
 - Compiler can rearrange instructions to make instruction-level parallelism by hardware more effective
 - If instruction-level parallelism in the instruction set (instructions with multiple operations in parallel) → compiler techniques to generate code for such machines from sequential programs
 - Parallelization techniques to translate sequential programs into *multiprocessor* code
- **Memory hierarchies:** registers + caches + physical memory + secondary storage
 - Cache-management by hw: not effective in scientific code operating on large arrays → compiler techniques for changing layout of data or order of instructions accessing data → improvement of effectiveness of memory hierarchy

Program Translations

- Translation between different kinds of languages
- **Binary Translation**: from binary code of one machine to binary code of another one (typically, to increase availability of software by computer companies)
- **Hardware Synthesis**: hw designs described in high-level hw description languages (VHDL, RTL) → translation by hw-synthesis tools into detailed hw schemas
- **Database Query Translators**: high-level query translated into actions for record search (SQL → Relational Algebra → Actions for physical search)
- **Compiled Simulation**: instead of writing a simulator interpreting the design (very expensive), better compiling the design into machine code that simulates the particular design

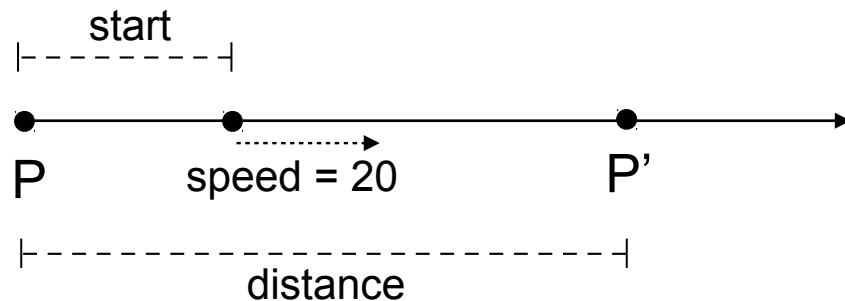
Context of a Compiler



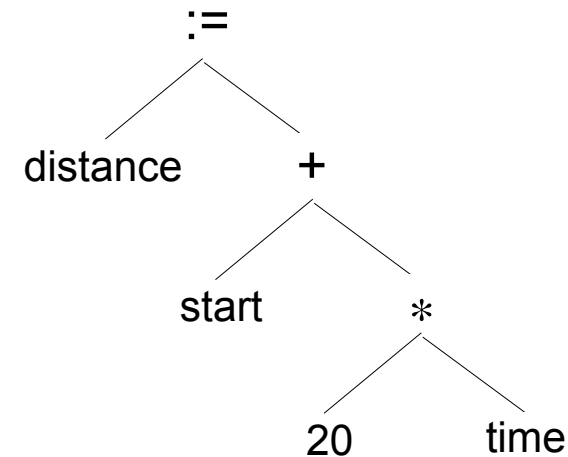
Compilation Model

- Separation < **analysis**: word recognition + internal representation of source (syntax tree)
synthesis: construction of target program → more specialized (and complex)


- Analysis → operations in source program mapped to hierarchical structure → **syntax tree**



```
distance := start + 20 * time
```

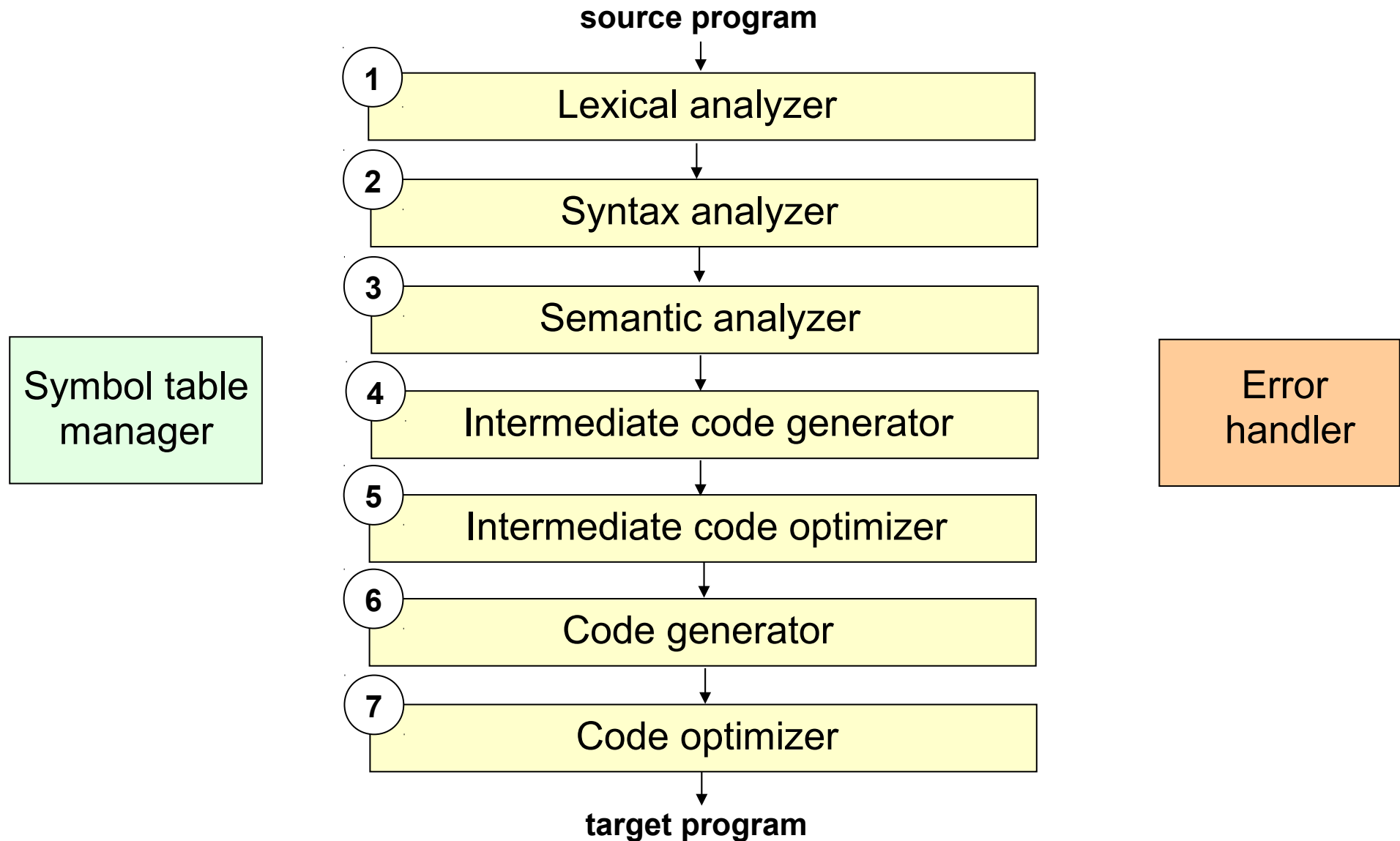


Analysis of Source Program

- **Lexical analysis** (linear) → grouping of characters into symbols → [token]
- **Syntax analysis** (hierarchical) → grammar symbols grouped hierarchically → syntax tree
- **Semantic analysis** → consistency checking 
 - type checking
 - insertion of information into symbol table

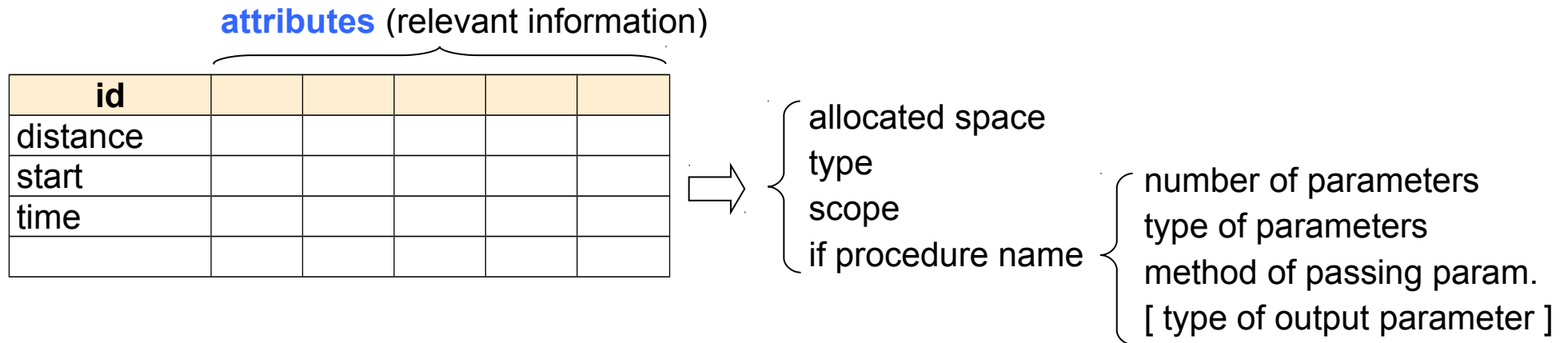
Phases of a Compiler

- **Phase** = conceptual unit in which a compiler operates: $\text{Rep}(\text{source } P) \rightarrow \text{Rep}'(\text{source } P)$



Symbol Table

- Data structure containing information on identifiers (catalog)



- Requirements {

efficient access to identifier attributes < read
 write

incremental update of attributes

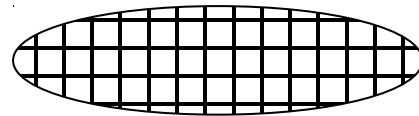
```
var distance, start, time: real;
```

Lexer → inserts identifiers, but
type known only later (sem)

- Example of use: < SEM: type checking
GEN: space allocated in memory

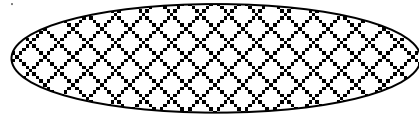
Error Handling

- \forall phase $F \rightarrow$ handling of errors pertinent to F (separation of concerns)



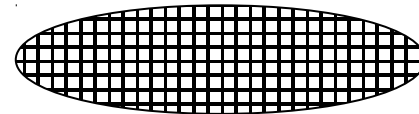
LEX

unknown symbol: @



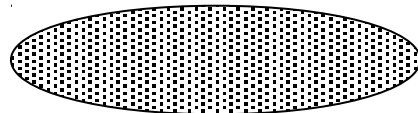
SYN

```
if a := 5 then
```



SEM

```
integer  
  ↙   ↘  
y := x + s ← string
```



RUNTIME

```
null  
  ↓  
*p := 10
```

Transformation of Source Program

- Same computation expressed by different abstraction levels

```
distance := start + 20 * time
```

1. Lexical analysis

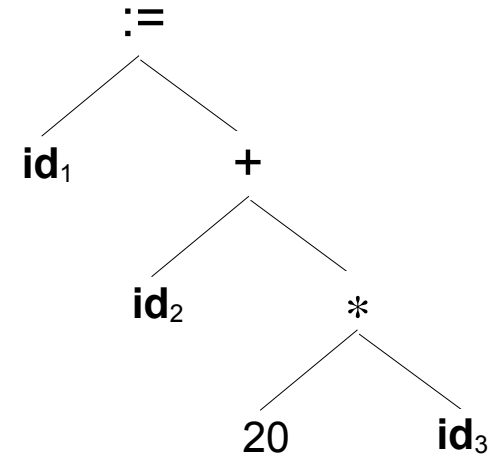
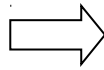
```
id1 := id2 + 20 * id3
```

- Recognition of (grammar) terminals → pairs (symbol, **value**)
- Removal of spacing / comments
- Encoding of each symbol (`:=` → `#define ASSIGN 257`)
- Some symbols enriched by a **lexical value** (**id** → “distance”, or pointer to symbol table)

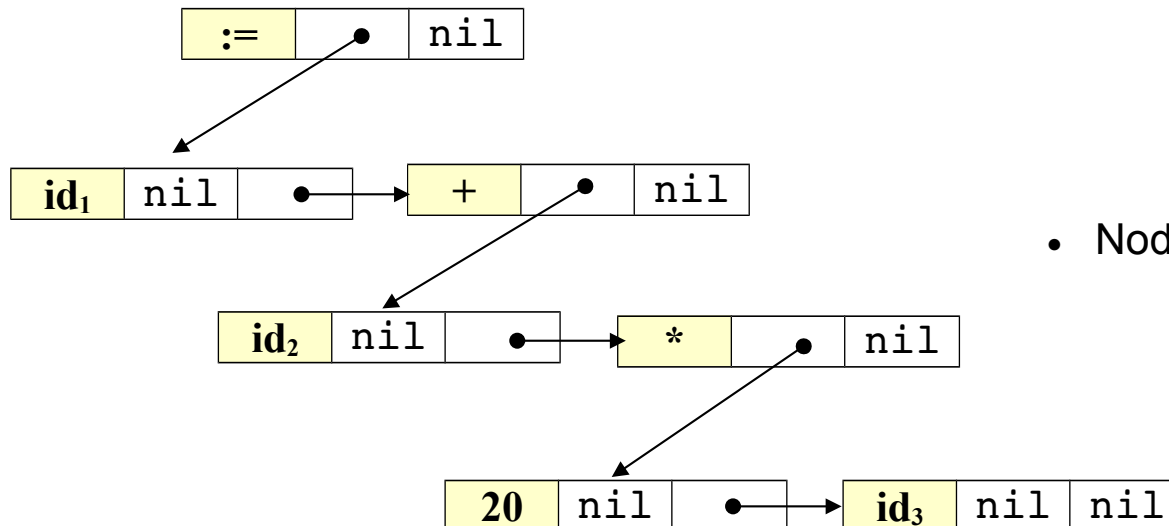
Transformation of Source Program (ii)

2. Syntax analysis

$id_1 := id_2 + 20 * id_3$



(token, child, brother)



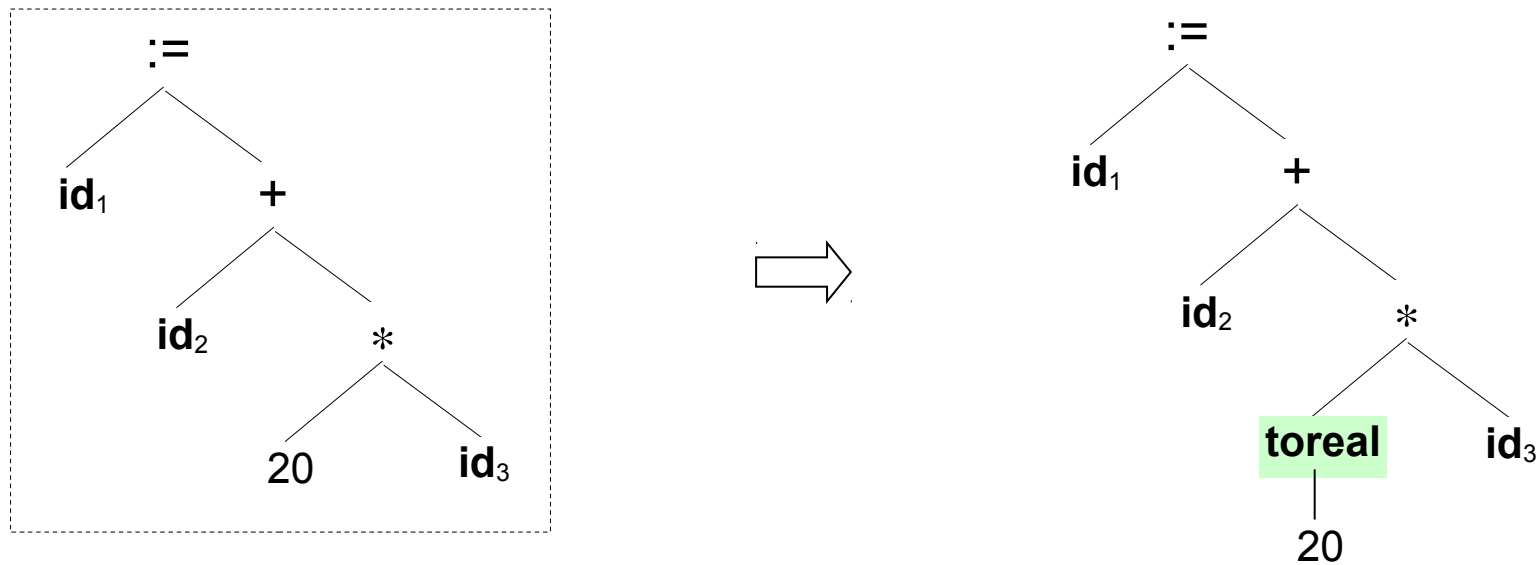
- Node uniformity → pb { inhomogeneity, child cardinality }

union

binarization

Transformation of Source Program (iii)

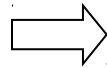
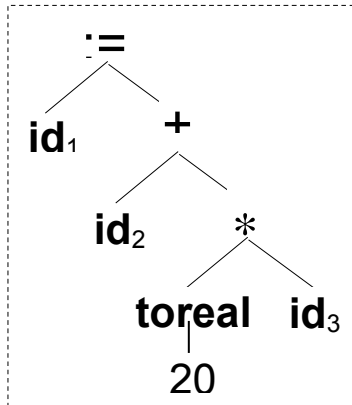
3. Semantic analysis



- Review of the tree for checking semantic constraints (types consistent with operations)
- Possible decoration/alteration of the tree

Transformation of Source Program (iv)

4. Intermediate code generation $\rightarrow \exists$ one statement for each operator of the syntax tree



```
t1 := toreal(20)
t2 := t1 * id3
t3 := id2 + t2
id1 := t3
```

\approx operational semantics

- Intermediate code = program written in a language of an abstract (virtual) machine
- Properties: easy to $\begin{cases} \text{generate} \\ \text{translate into target code} \end{cases}$
- Advantages $\begin{cases} \text{implementation simplification} \\ \text{portability (reusability)} \end{cases}$
- Nature: varying, typically: three-address code (quadruples) \approx Assembly where memory locations viewed as registers

operator addr₁ addr₂ addr₃



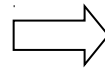
* addr₁ addr₂ addr₃ = addr₃ := addr₁ * addr₂

- \exists at most one explicit operator (in addition to assignment) \rightarrow operation linearization
- Generation of temporaries for intermediate results
- Not necessarily three operands (possibly less, e.g. **toreal**)

Transformation of Source Program (v)

5. Intermediate code optimization → reduction of number of instructions

```
t1 := toreal(20)
t2 := t1 * id3
t3 := id2 + t2
id1 := t3
```



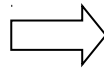
```
t := 20.0 * id3
id1 := id2 + t
```

- Intermediate code more efficient (e.g.: conversion 20 → 20.0 can be performed statically!)
- Problem: slowdown of compiling (for advanced optimizations)
- Possible different optimization levels
- Better not to optimize during coding (time, debugging)
- After optimization → redo testing

Transformation of Source Program (vi)

6. **Code generation** → need for register-loading to perform operations

```
t := 20.0 * id3
id1 := id2 + t
```

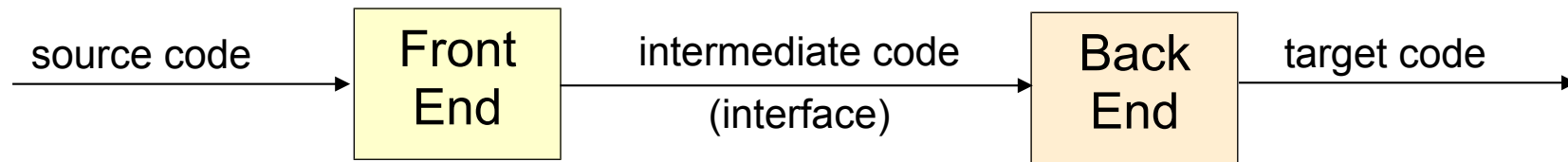


LDF	R ₂ ,	id ₃	
MULF	R ₂ ,	R ₂ ,	#20.0
LDF	R ₁ ,	id ₂	
ADDF	R ₁ ,	R ₁ ,	R ₂
STF	id ₁ ,	R ₁	

- Type of target code
 - Assembly (typically)
 - Machine, relocatable → variables transformed into memory locations
- In general: mapping
 - statements → equivalent machine statements
 - variables → registers (for operations)

Grouping of Compiler Phases

- Logical organization \neq physical organization (like in DB)



- **Front-End**: depends only on source \rightarrow LEX, SYN, SEM, GEN-I [, part of OPT-I]
- **Back-End**: depends only on target \rightarrow OPT-I, GEN, OPT
- Macro-modularization: useful for the porting of compiler (reusable macro-modules)

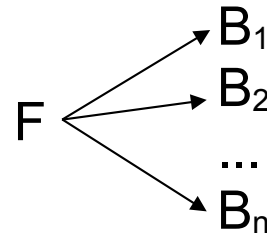


Ideally: changing $\left\{ \begin{array}{l} \text{Source L} \rightarrow \text{change of Front-End} \\ \text{Target L} \rightarrow \text{change of Back-End} \end{array} \right.$ \Longleftarrow intermediate code: invariant

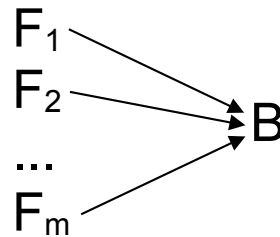
Grouping of Compiler Phases (ii)

- Possible scenarios:

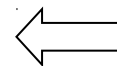
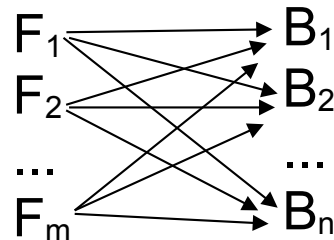
1. Same L on \neq platforms (realistic):



2. \neq L on same platform (ideal):



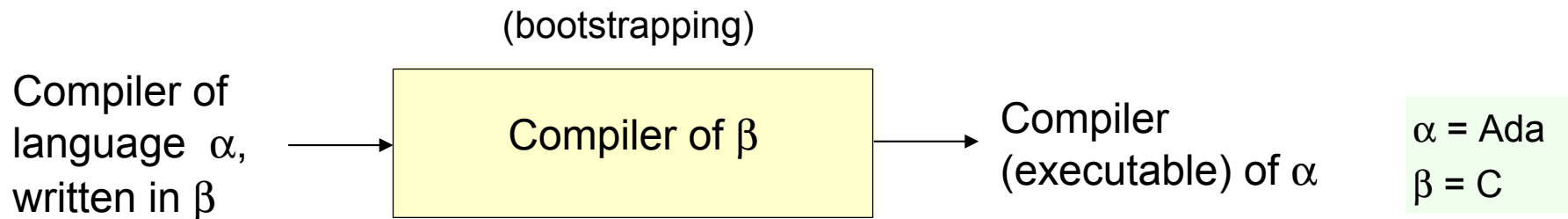
3. \neq L on \neq platforms:



$(n+m)$ modules rather than $(m*n)$ compilers

Bootstrapping & Porting

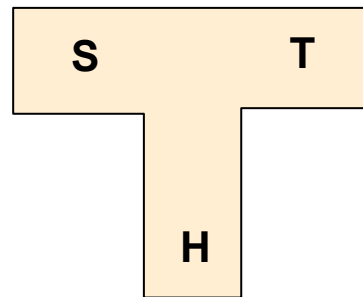
- L involved in compiler construction
 - source
 - implementation (host)
 - target
- Machine language in first compilers*



- **Cross-compiler**: when $\text{Comp}(\alpha)$ runs on a machine \neq target
(necessary when target machine with limited resources)
→ target code not executable on the same machine
- target = Intel
 $\text{Comp}(\alpha)$ on AMD

Bootstrapping & Porting (ii)

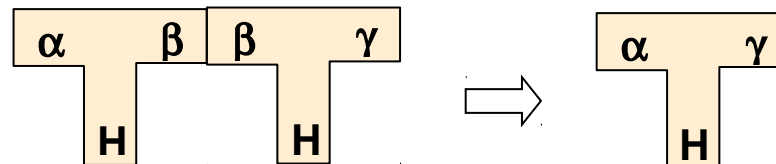
- **T-diagram**: to schematize a compiler (by means of the 3 involved languages)



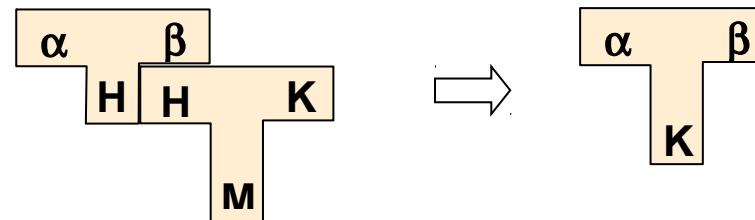
- \underline{H} = Machine language of executable
- $\underline{H} \neq T \rightarrow$ cross-compiler

- Diagram composition (to generate new compilers):

1. **Concatenation**
(Change of *target* language)



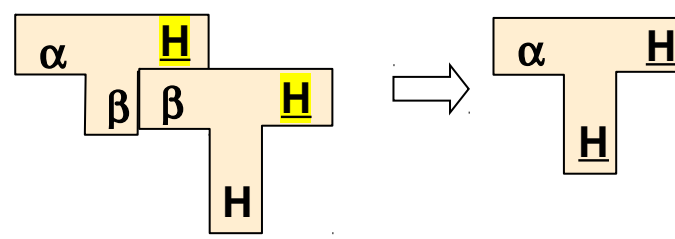
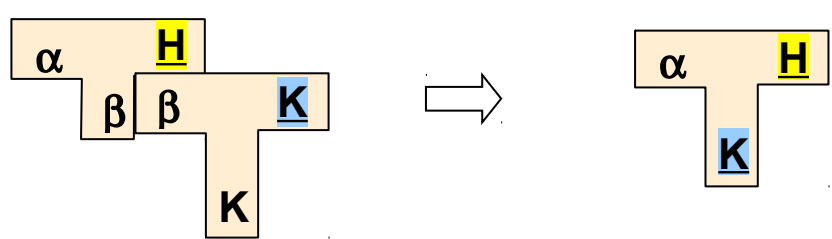
2. **Joint**
(Change of *host* language)

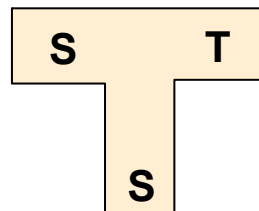


α = Ada
 β = C
H = Pascal
K = Fortran

Bootstrapping & Porting (iii)

Notes (3 kinds of bootstrapping):

- Compiler of language α , written in β \rightarrow (bootstrapping)
Compiler of β \rightarrow Compiler (executable) of α \Rightarrow 
- Bootstrapping of cross-compiler \Rightarrow 
- Compiler written in the source language to compile ! \rightarrow Pb of bootstrapping !**



Bootstrapping & Porting (iv)

- Solution of the circularity: $\left\{ \begin{array}{l} \text{GOOD} = \text{compiler of } S \text{ written in } S \begin{cases} \text{efficient} \\ \text{generates efficient target code} \end{cases} \\ \underline{\text{BAD}} = \text{compiler of } S' \subset S \text{ written in Assembly} \end{array} \right.$
- Limits of BAD: $\left\{ \begin{array}{l} \text{Compiles only a subset of } S \text{ (relevant to implementation of GOOD)} \\ \text{Inefficient at runtime} \\ \text{Generates inefficient target code} \begin{cases} \text{space} \\ \text{time} \end{cases} \end{array} \right.$
- Unique requirement for BAD: *functional correctness of the generated code*



Bootstrapping & Porting (v)

Bootstrapping (alchemy):

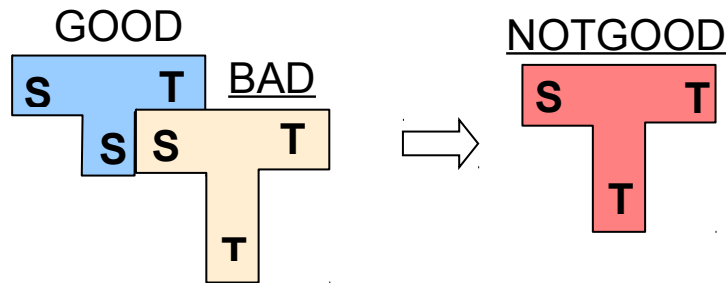
1. Writing in Assembly of a compiler “*quick & dirty*” (BAD) for a subset of S used to write GOOD

2. Compiling of GOOD using BAD → NOTGOOD

inefficient (runtime)

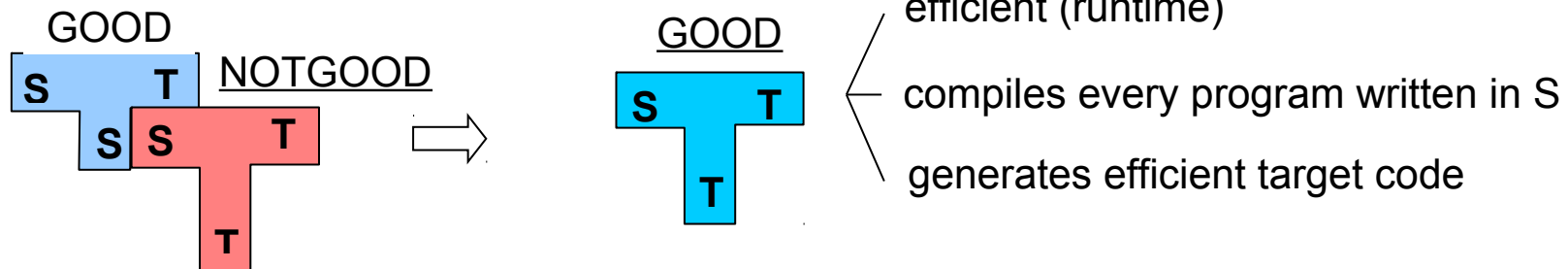
compiles every program written in S

generates efficient target code!



Requirement of functional correctness for the code generated by BAD

3. Compiling of GOOD using NOTGOOD → GOOD

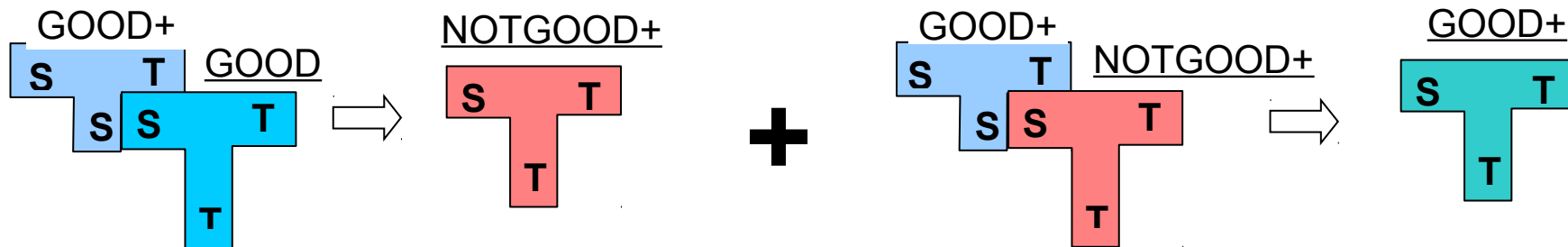


Bootstrapping & Porting (vi)

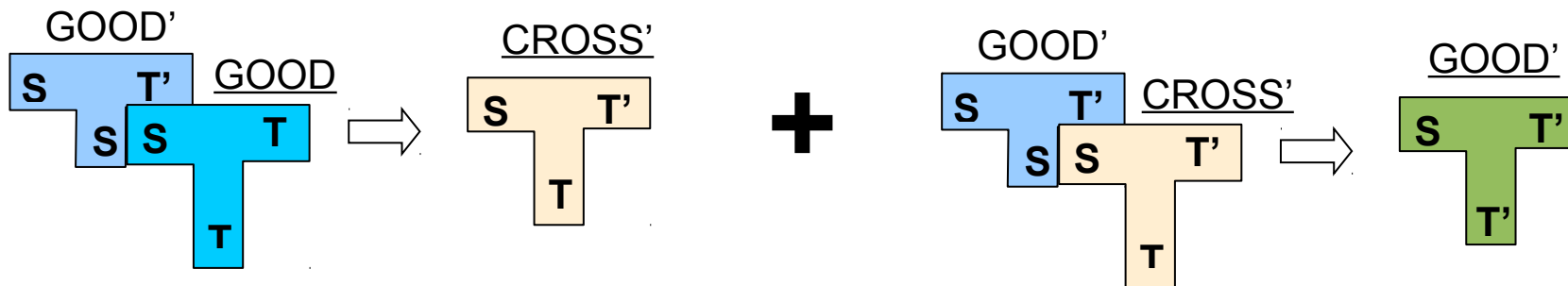
- After bootstrapping → compiler specified in 2 forms $\begin{cases} \text{source } S = \text{GOOD} \\ \text{target } T = \underline{\text{GOOD}} \end{cases}$

Advantages:

- a) Improvement of GOOD → GOOD+ → immediately bootstrapped by steps 2 and 3 in GOOD+



- b) Porting of GOOD to a new host computer → only rewritten the Back-End of GOOD → GOOD'



Bootstrapping & Porting (vii)

- Similar considerations when “*quick & dirty*” written in a high-level language \neq S

