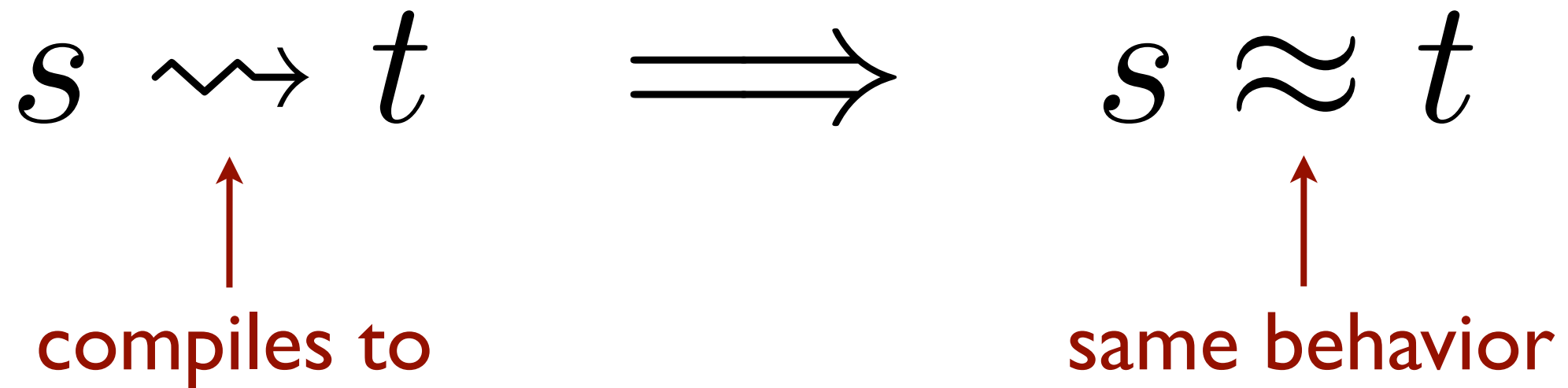


Correct and Secure Compilation for Multi-Language Software

Amal Ahmed

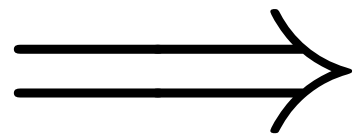
Northeastern University

Compiler Correctness



Semantics-preserving compilation

$s \rightsquigarrow t$
↑
compiles to



$s \approx t$
↑
same behavior

Range of Compiler Properties...

- Type-preserving compilation (90s)
- Semantics-preserving compilation (00s...)
= *Correct compilation*
- Fully abstract compilation
= Equivalence-preserving and -reflecting
= *Secure compilation*
- Security-preserving compilation
 - preserving “security types” vs. preserving noninterference

Compiler Verification

One of the “big problems” of computer science

- since *McCarthy and Painter 1967*:
Correctness of a Compiler for Arithmetic Expressions
- see *Dave 2003: Compiler Verification: A Bibliography*

Compiler Verification since 2006...

*Leroy '06 : Formal certification of a compiler back-end or:
programming a compiler with a proof assistant.*

Lochbihler '10 : Verifying a compiler for Java threads.

Myreen '10 : Verified just-in-time compiler on x86.

*Sevcik et al.'11 : Relaxed-memory concurrency and verified
compilation.*

*Zhao et al.'13 : Formal verification of SSA-based
optimizations for LLVM*

Kumar et al.'14 : CakeML: A verified implementation of ML

⋮

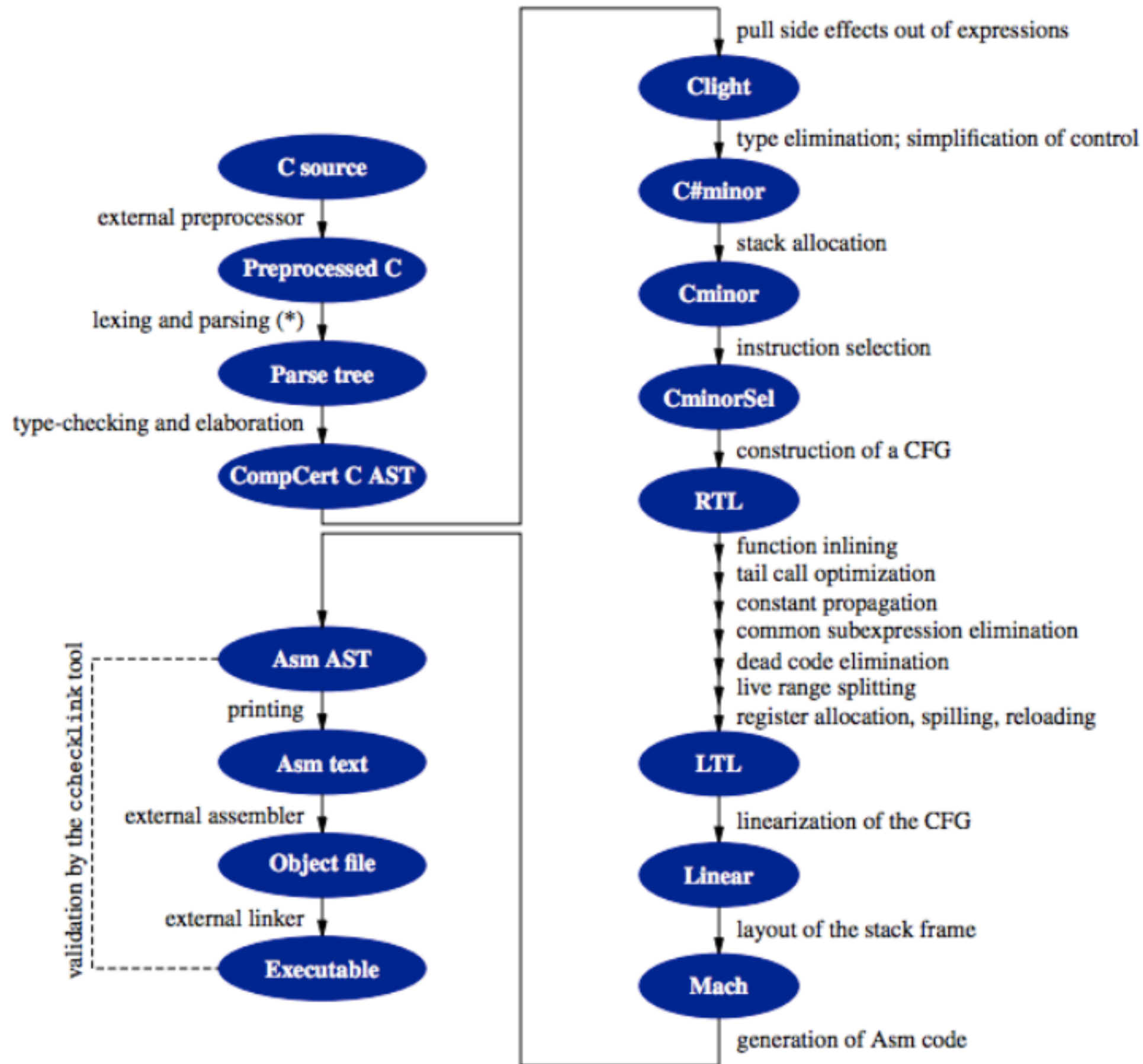
Why CompCert had such impact...

- Demonstrated that realistic verified compilers are both *feasible* and bring *tangible benefits*

The striking thing about our CompCert results is that the middle-end bugs we found in all other compilers are absent. As of early 2011, the under-development version of CompCert is the only compiler we have tested for which Csmith cannot find wrong-code errors. This is not for lack of trying: we have devoted about six CPU-years to the task. The apparent unbreakability of CompCert supports a strong argument that developing compiler optimizations within a proof framework, where safety checks are explicit and machine-checked, has tangible benefits for compiler users. (Yang et al. PLDI 2011)

Why CompCert had such impact...

- Provided a proof architecture for others to follow/build on
 - CompCert memory model, uniform across passes
 - proof using simulations



Not verified yet
 (*) the parser is formally verified

Formally verified

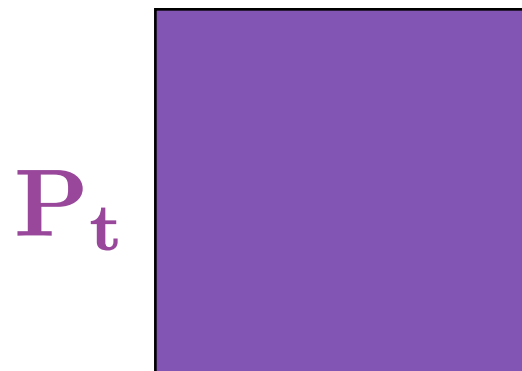
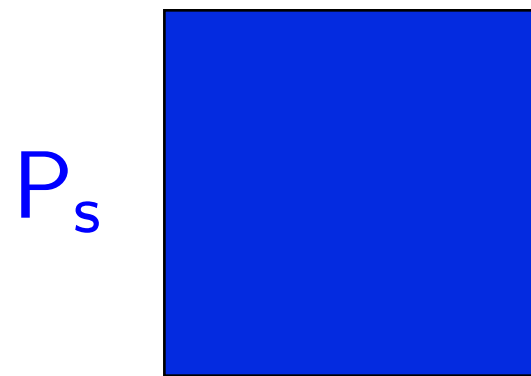
Why CompCert had such impact...

- Provided a proof architecture for others to follow/build on
 - CompCert memory model, uniform across passes
 - proof using simulations

*But the simplicity of the proof architecture
comes at a price...*

Problem: Whole-Program Assumption

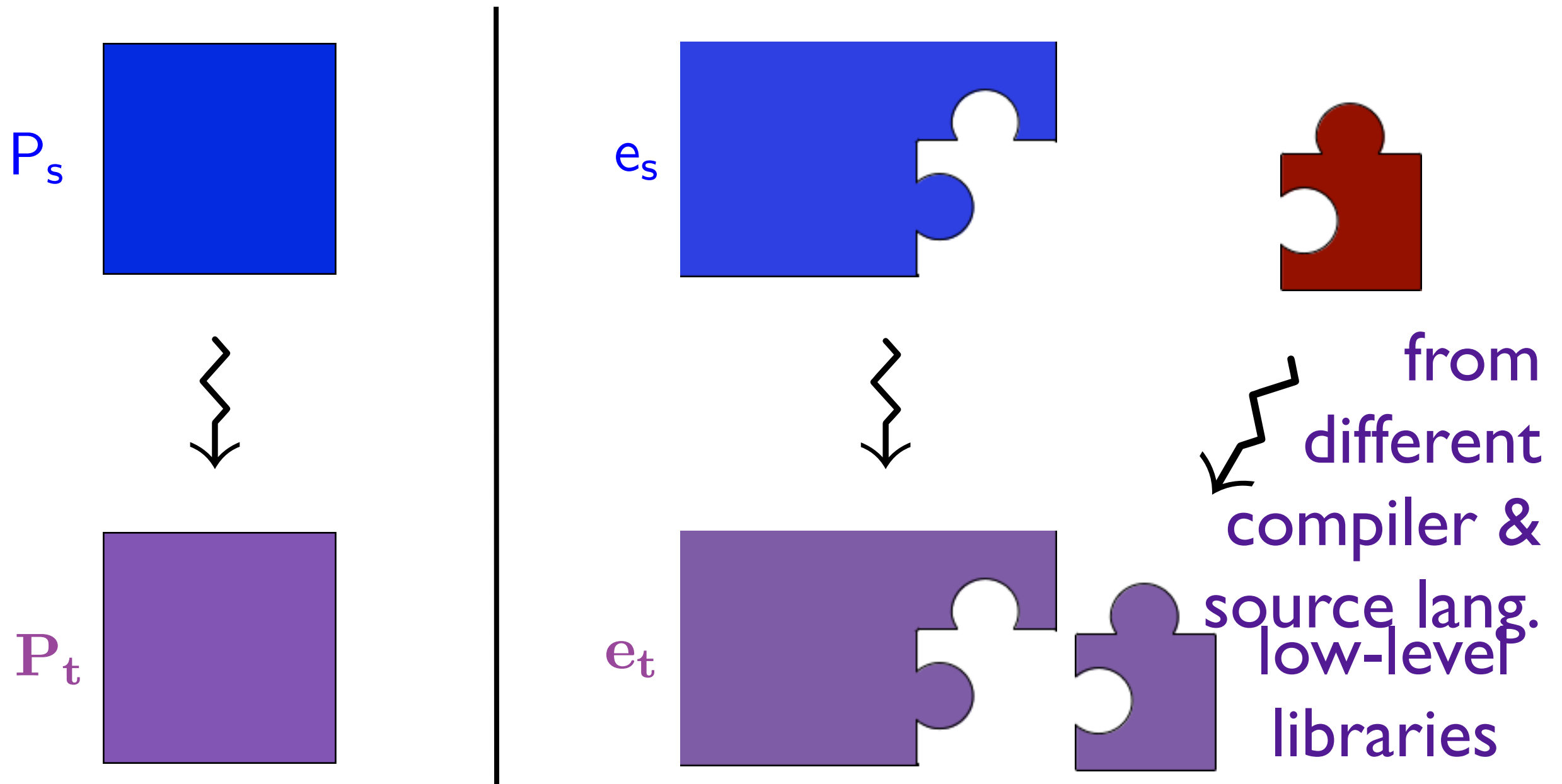
Correct compilation guarantee only applies to **whole** programs!



CompCert's ... “formal guarantees of semantics preservation apply only to whole programs that have been compiled as a whole by [the] CompCert C [compiler]” (Leroy 2014)

Problem: Whole-Program Assumption

Correct compilation guarantee only applies to **whole** programs!



Why Whole Programs?

$$s \rightsquigarrow t \quad \Longrightarrow \quad s \approx t$$

↑
expressed how?

Why Whole Programs?

$$P_s \rightsquigarrow P_t \implies P_s \approx P_t$$

↑
expressed how?

CompCert

$$\begin{array}{ccccccc} P_s & \mapsto & \dots & \mapsto & P_s^i & \mapsto & P_s^{i+1} & \mapsto & \dots \\ \vdots & & & & \vdots & & \vdots & & \\ P_t & \mapsto & \dots & \mapsto & P_t^j & \mapsto^* & P_t^{j+n} & \mapsto & \dots \end{array}$$

Proof composes per-pass simulations

$$P_s \rightsquigarrow P_t \implies \begin{array}{c} P_s \longmapsto \dots \longmapsto P_s^i \longmapsto P_s^{i+1} \longmapsto \dots \\ \text{\color{red} | } R_1 \\ P_t \longmapsto \dots \longmapsto P_t^j \longmapsto^* P_t^{j+n} \longmapsto \dots \end{array}$$

$$P_t \rightsquigarrow P_u \implies \begin{array}{c} P_t \longmapsto \dots \longmapsto P_t^i \longmapsto P_t^{i+1} \longmapsto \dots \\ \text{\color{red} | } R_2 \\ P_u \longmapsto \dots \longmapsto P_u^j \longmapsto^* P_u^{j+n} \longmapsto \dots \end{array}$$

$$P_s \rightsquigarrow P_u \implies P_s \approx P_u$$

Why Whole Programs?

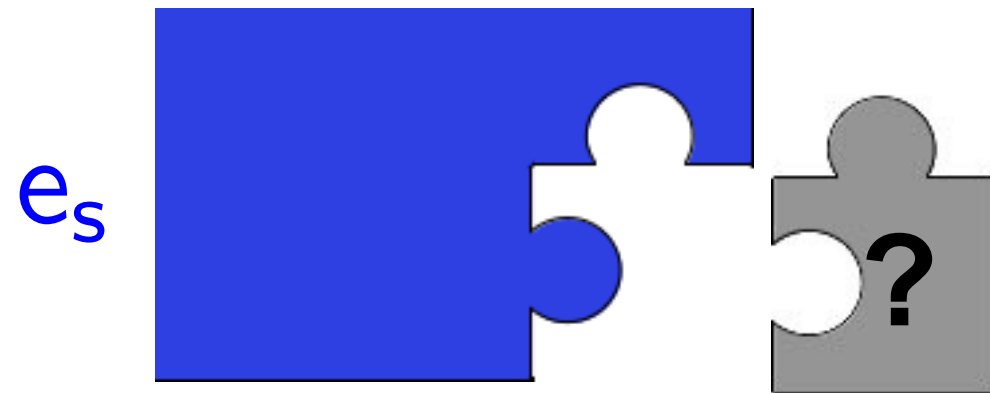
$$P_s \rightsquigarrow P_t \implies P_s \approx P_t$$

↑
“closed” simulations

CompCert

$$\begin{array}{ccccccc} P_s & \mapsto & \dots & \mapsto & P_s^i & \mapsto & P_s^{i+1} & \mapsto & \dots \\ \vdots & & & & \vdots & & \vdots & & \\ P_t & \mapsto & \dots & \mapsto & P_t^j & \mapsto^* & P_t^{j+n} & \mapsto & \dots \end{array}$$

Correct Compilation of Components?



e'_t

$$e_s \approx e_T$$



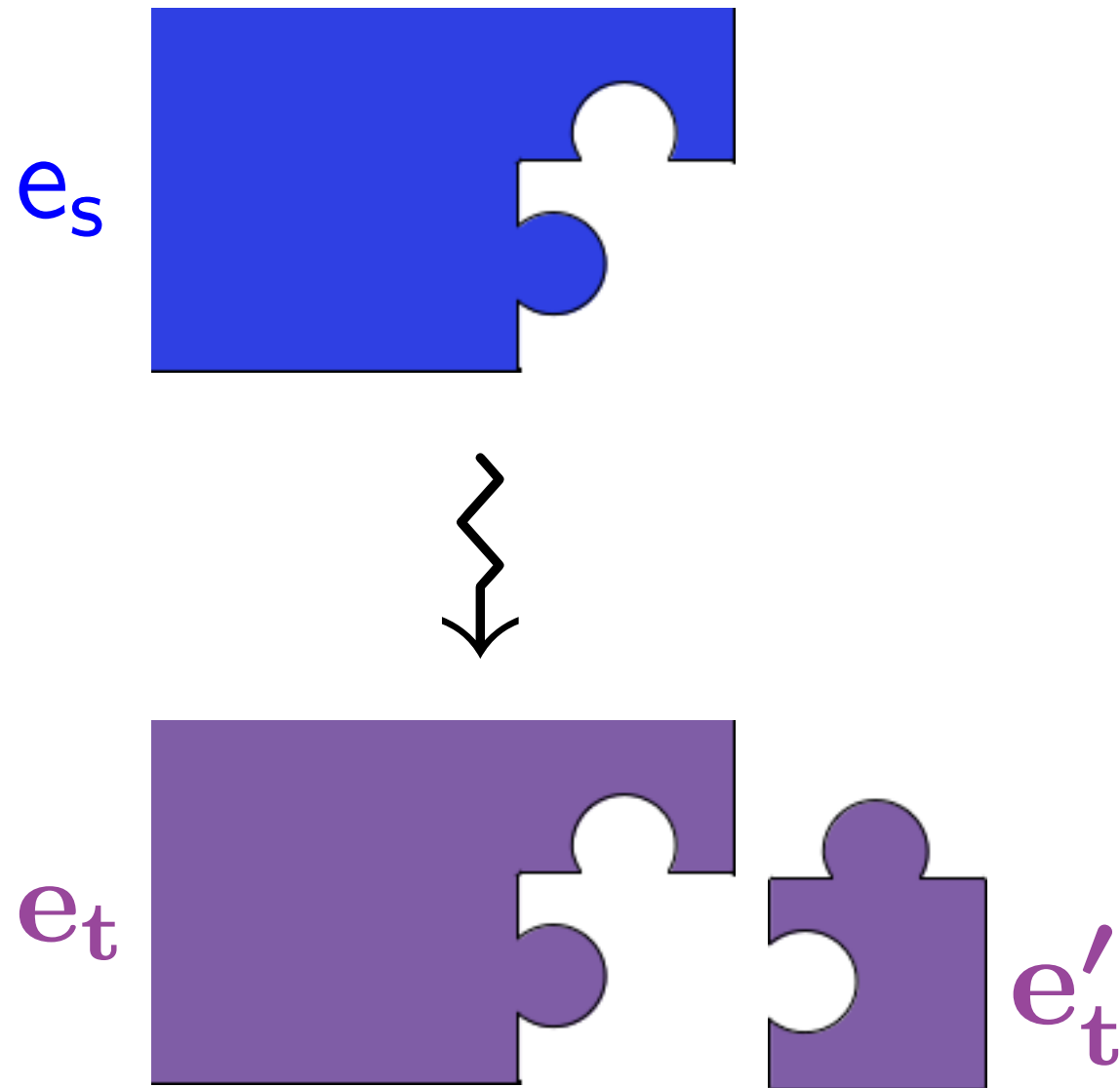
expressed how?

Produced by

- same compiler,
- diff compiler for S ,
- compiler for diff lang R ,
- R that's **very** diff from S ?

Behavior expressible in S ?

Correct Compilers, Multi-language SW



$$e_s \approx e_T$$



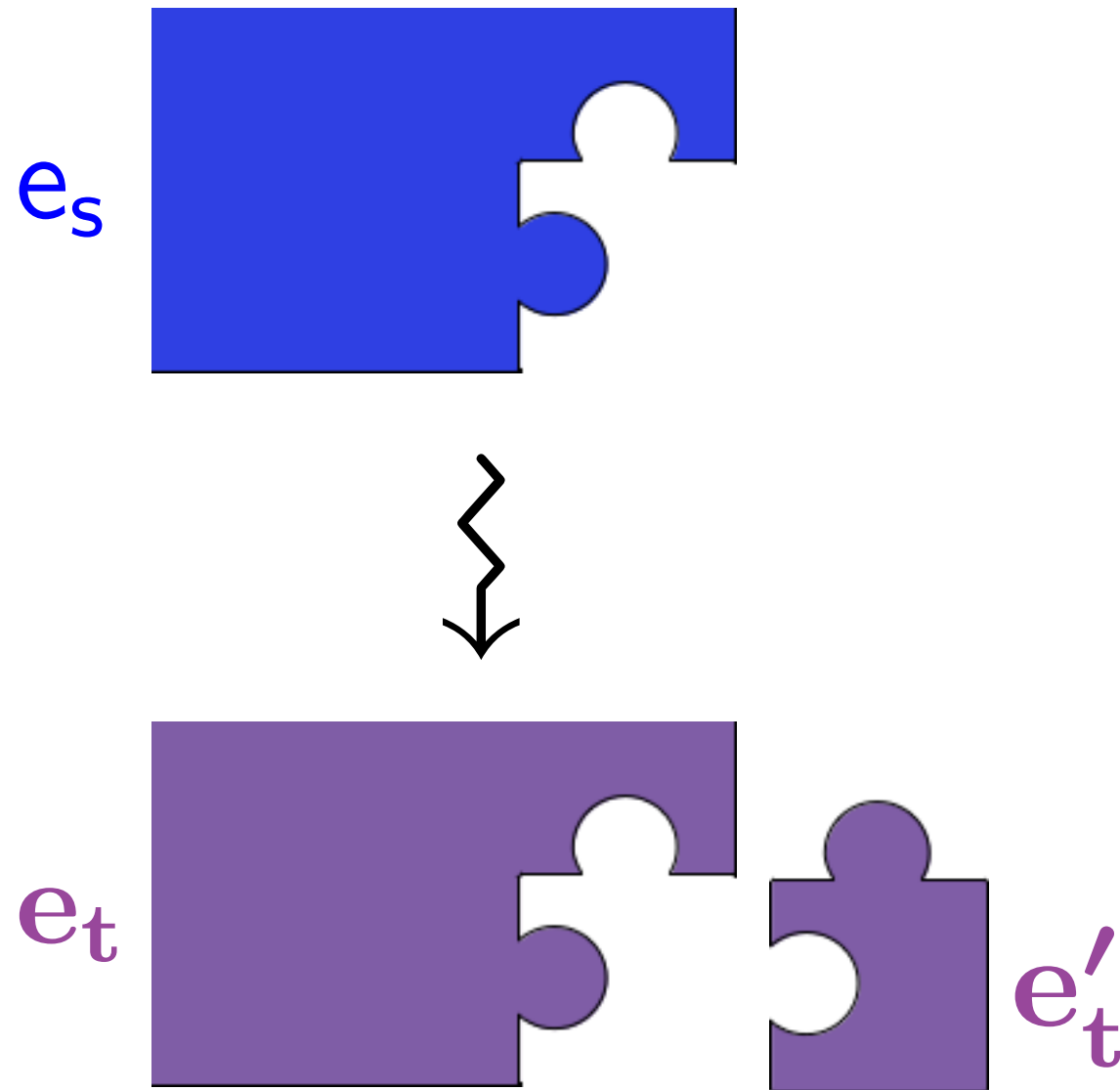
Definition should:

- permit **linking** with target code of arbitrary provenance
- support verification of **multi-pass** compilers

Plan

- Survey the literature: how to express $e_S \approx e_T$
 - “compositional” compiler correctness
 - = correct compilation of components

Compositional Compiler Correctness



$$e_s \approx e_T$$

Dictates:

- what we can **link** with (*horizontal compositionality*) and how to check it's okay to link
- effort involved in proving *transitivity* for **multi-pass** compilers (*vertical compositionality*)

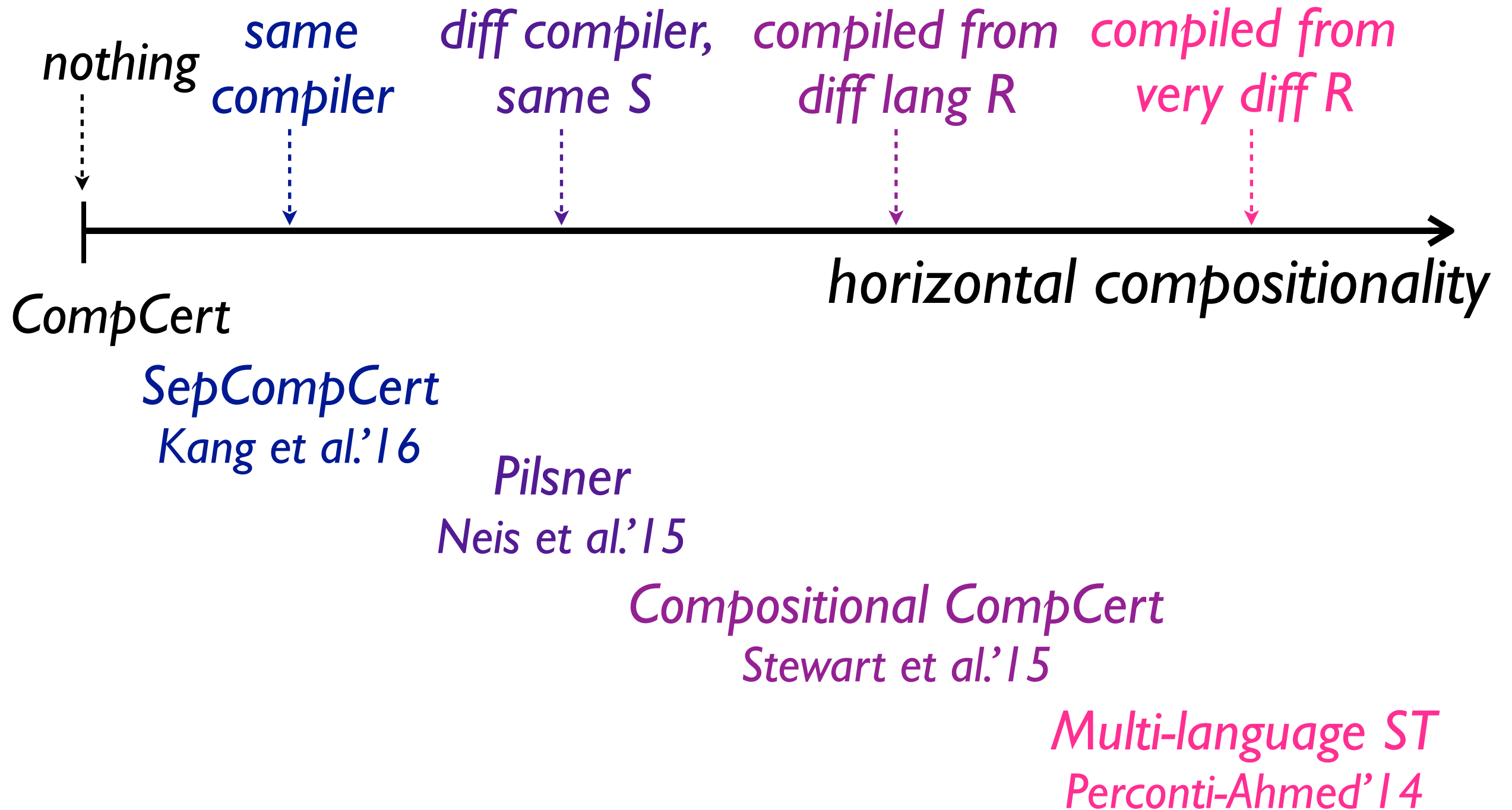
Plan

- Survey the literature: how to express $e_S \approx e_T$
- How does the choice affect:
 - what we can link with
 - how we check if some e'_t is okay to link with
 - effort required to prove transitivity
 - effort required to have confidence in theorem statement
- How to support linking with code from **very** different R

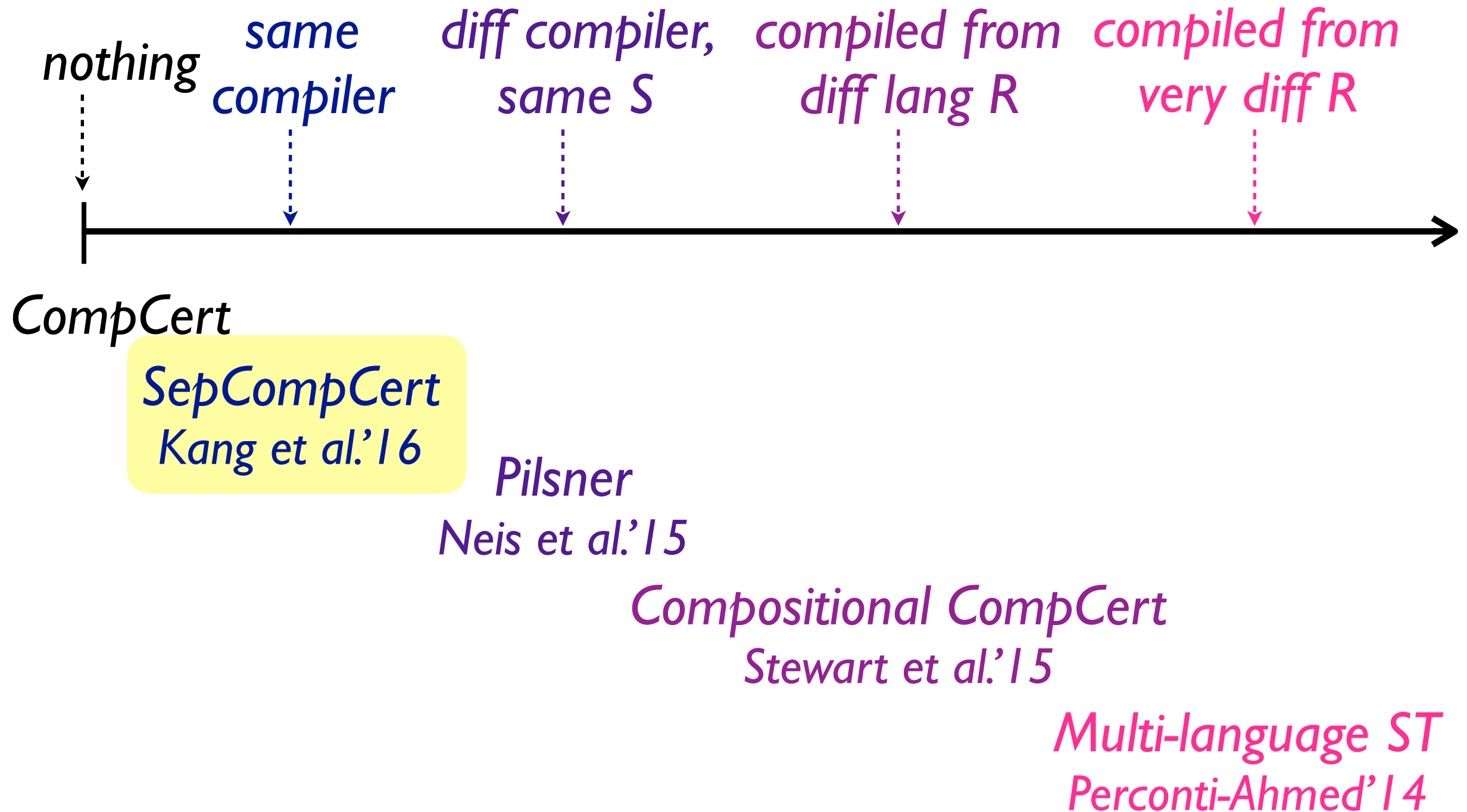
Plan

- Survey the literature: how to express $e_S \approx e_T$
- How does the choice affect:
 - what we can link with
 - how we check if some e'_t is okay to link with
 - effort required to prove transitivity
 - effort required to have confidence in theorem statement
- How to support linking with code from **very** different R
- Type-preserving compilation
- Secure (fully abstract) compilation

What we can link with



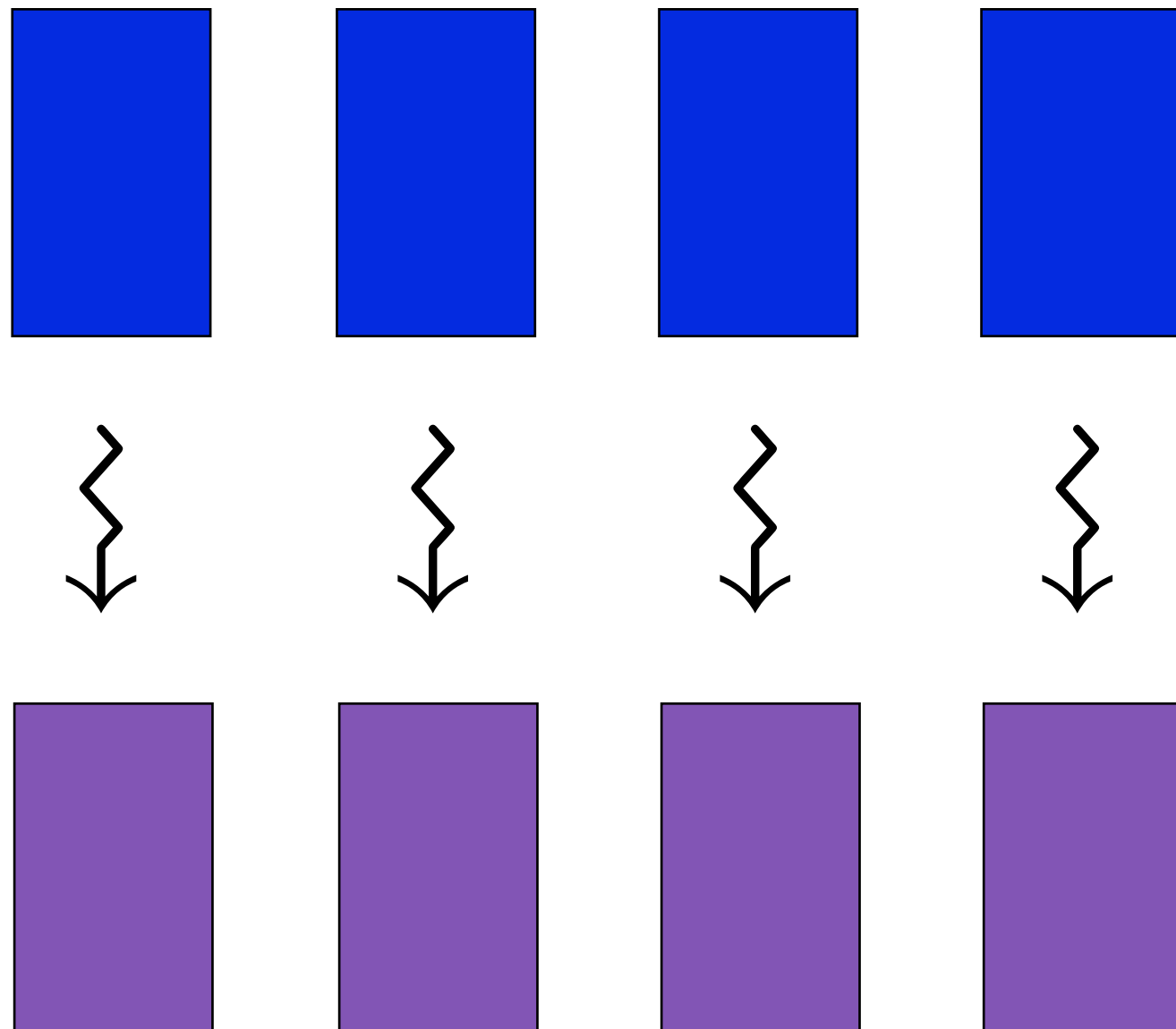
What we can link with



Approach: Separate Compilation (C)

SepCompCert

[Kang et al. '16]



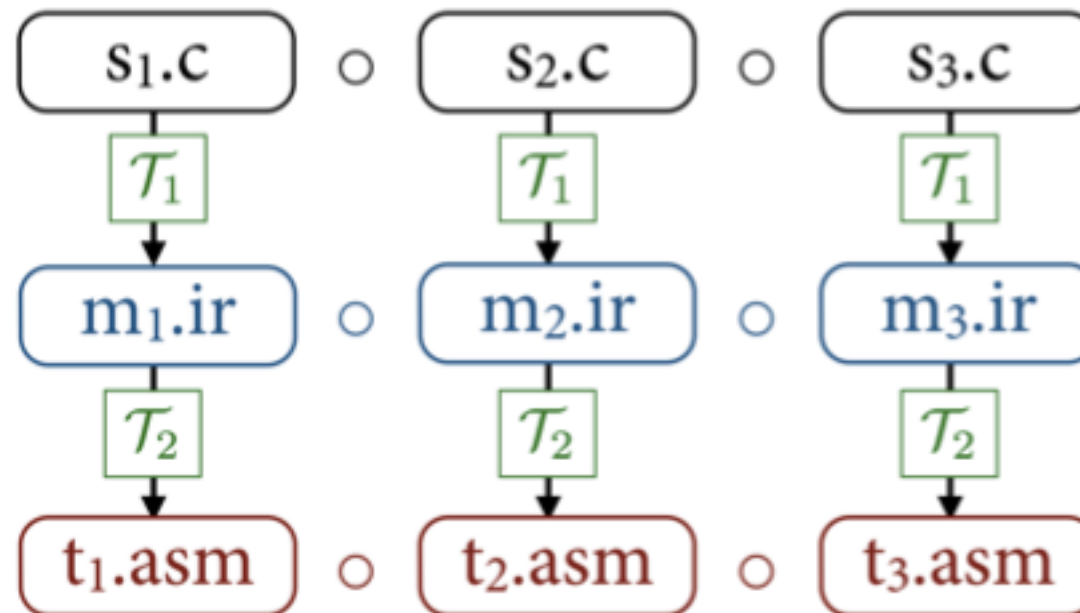
Approach: Separate Compilation (C)

SepCompCert

[Kang et al. '16]

Level A correctness

End-to-end

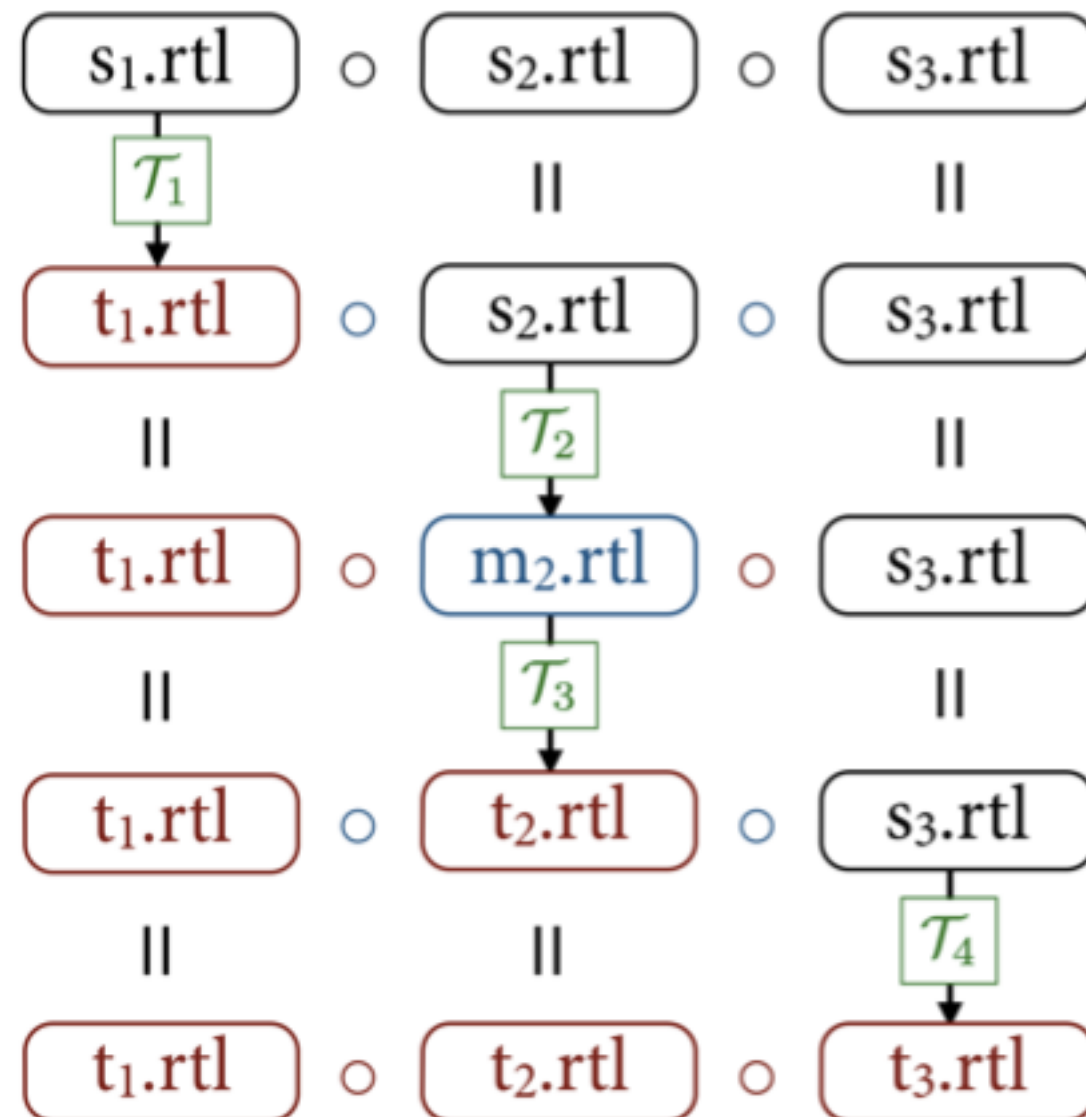


$$\frac{\forall i \in \{1 \dots n\}. \mathcal{C}(s_i.c) = t_i.asm \quad s = \text{load}(s_1.c \circ \dots \circ s_n.c) \quad t = \text{load}(t_1.asm \circ \dots \circ t_n.asm)}{\text{Behav}(s) \supseteq \text{Behav}(t)}$$

Approach: Separate Compilation (C)

SepCompCert
[Kang et al. '16]

Level B correctness: omit some RTL optimizations



Approach: Separate Compilation (C)

SepCompCert

[Kang et al. '16]

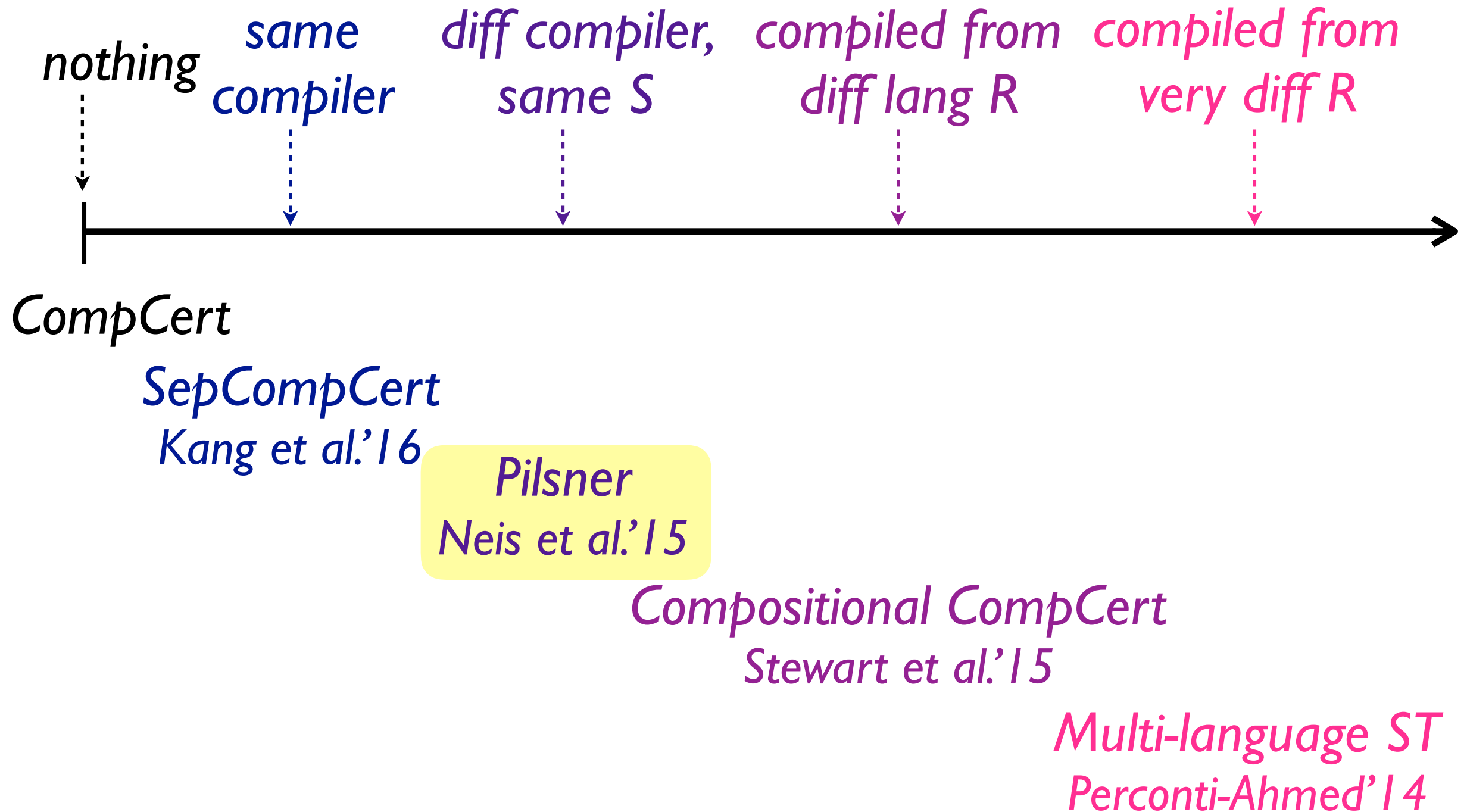
Level B correctness: omit some RTL
optimizations

End-to-end

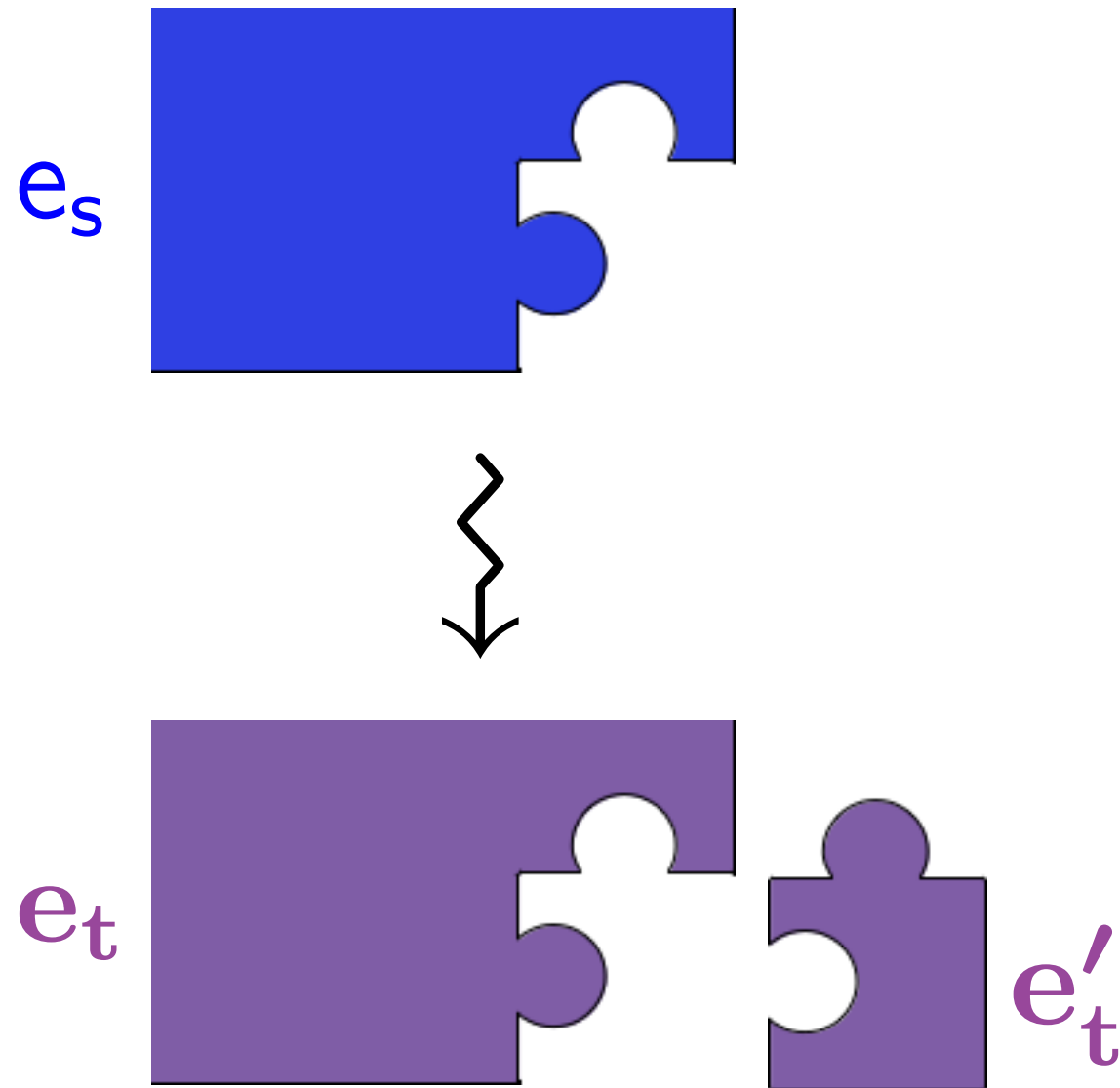


$$\frac{\forall i \in \{1 \dots n\}. C_i(\mathbf{s}_i.c) = \mathbf{t}_i.asm}{s = \text{load}(\mathbf{s}_1.c \circ \dots \circ \mathbf{s}_n.c) \quad t = \text{load}(\mathbf{t}_1.asm \circ \dots \circ \mathbf{t}_n.asm)} \text{Behav}(s) \supseteq \text{Behav}(t)$$

What we can link with



Approach: Cross-Language Relations



Cross-language relation

$$e_s \approx e_T$$

Compiling ML-like langs:

Logical relations

- [Benton-Hur ICFP'09]
- [Hur-Dreyer POPL'11]

Parametric inter-language simulations (PILS)

- [Neis et al. ICFP'15]

Case Study: Closure Conversion

- Typed Closure Conversion
- Correctness of closure conversion
using a cross-language logical relation...

[on board]

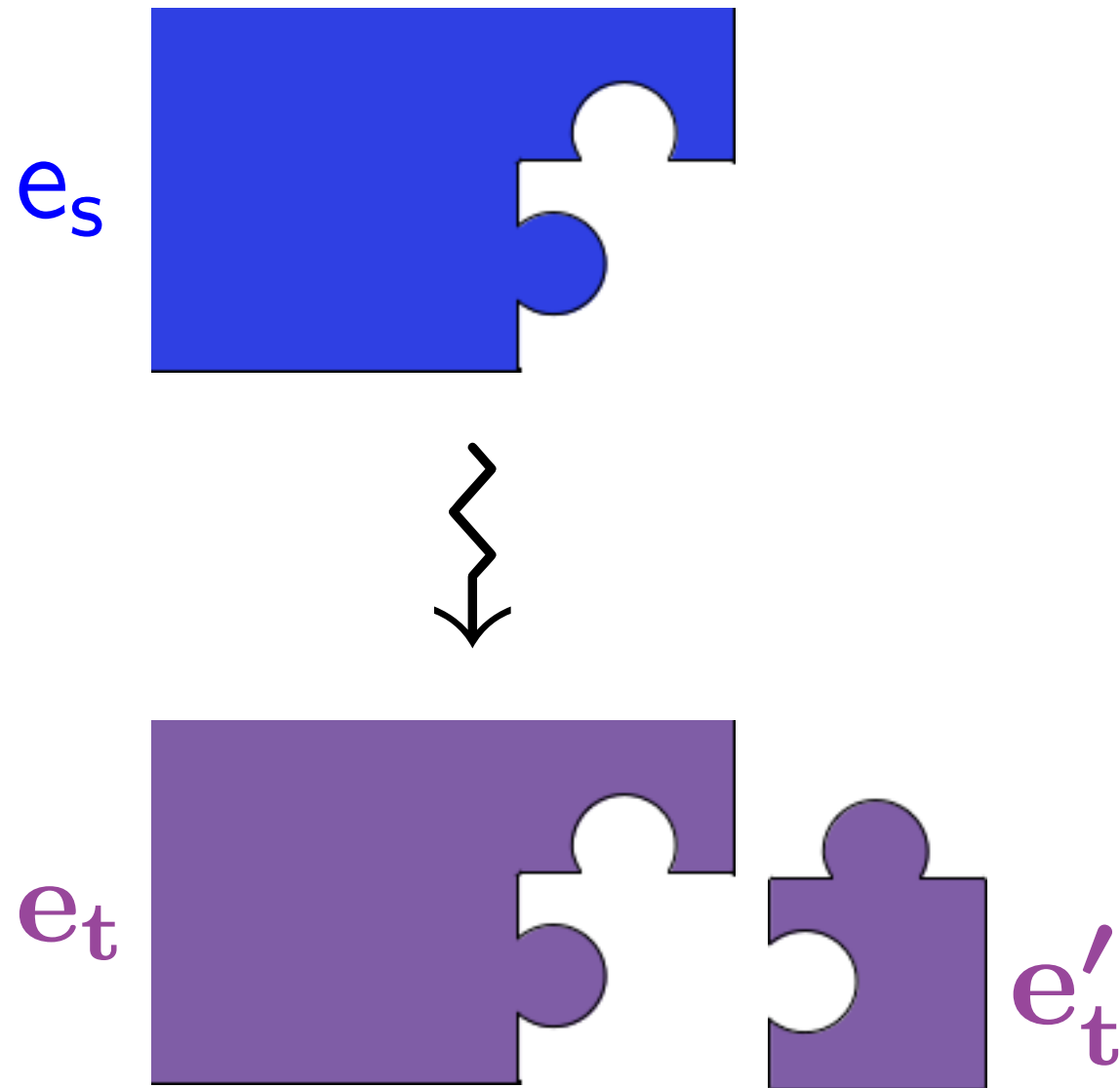
Cross-Language Relation: Problem 1

$$x : \tau' \vdash e_s : \tau \rightsquigarrow e_t \implies x : \tau' \vdash e_s \simeq e_t : \tau$$

cross-language logical relation

$$\forall e'_s, e'_t. \vdash e'_s \simeq e'_t : \tau' \implies \vdash e_s[e'_s/x] \simeq e_t[e'_t/x] : \tau$$

Cross-Language Relation: Problem 1



Have $x : \tau' \vdash e_s \simeq e_t : \tau$

*Does the compiler
correctness theorem
permit linking with e'_t ?*

Cross-Language Relation: Problem 1

Have $x : \tau' \vdash e_s \simeq e_t : \tau$



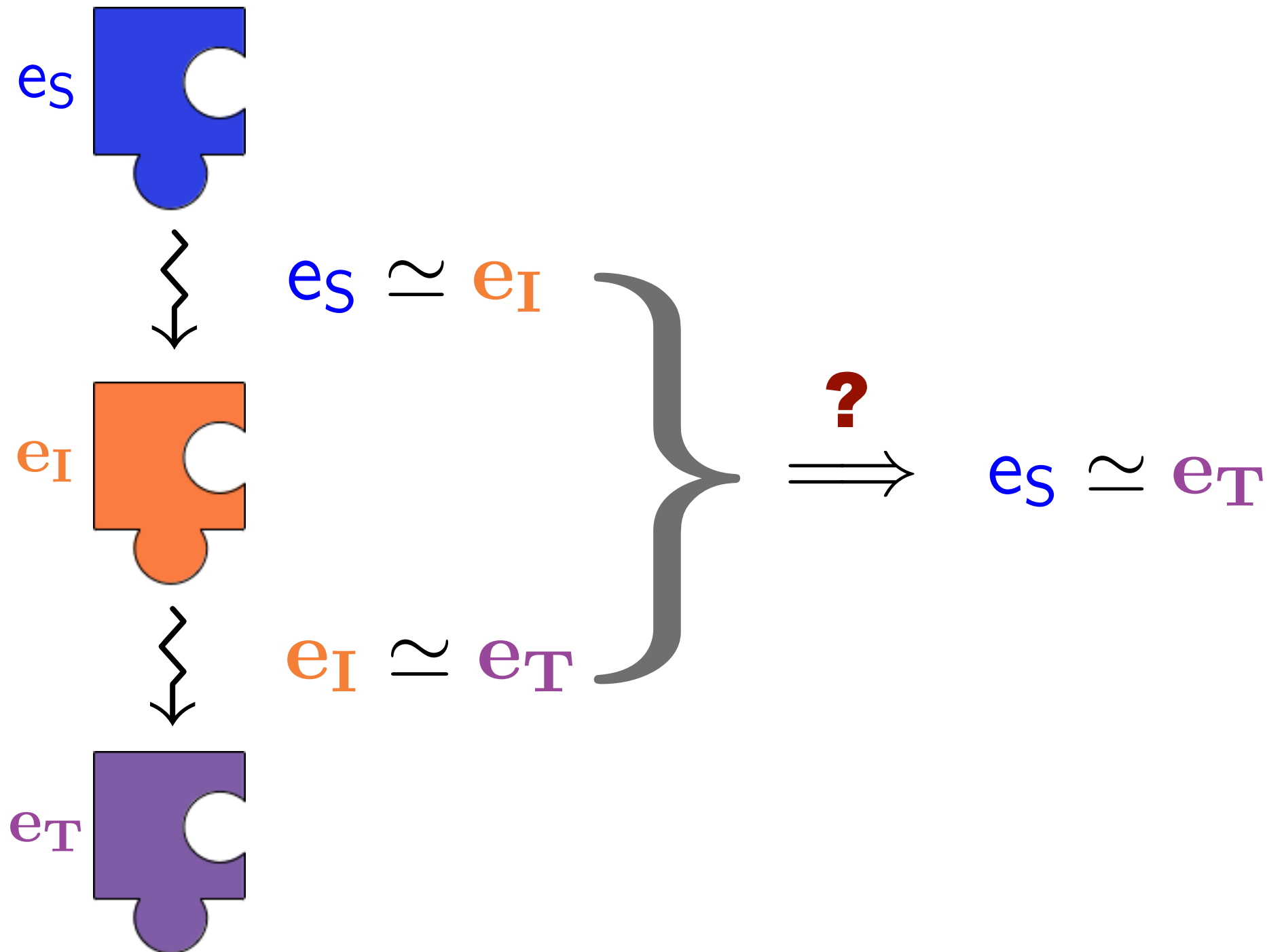
cross-language logical relation

$$\forall e'_s, e'_t. \vdash e'_s \simeq e'_t : \tau' \implies \vdash e_s[e'_s/x] \simeq e_t[e'_t/x] : \tau$$



- Need to come up with e'_s
-- not feasible in practice!
- Cannot link with e'_t
whose behavior cannot be expressed in source.

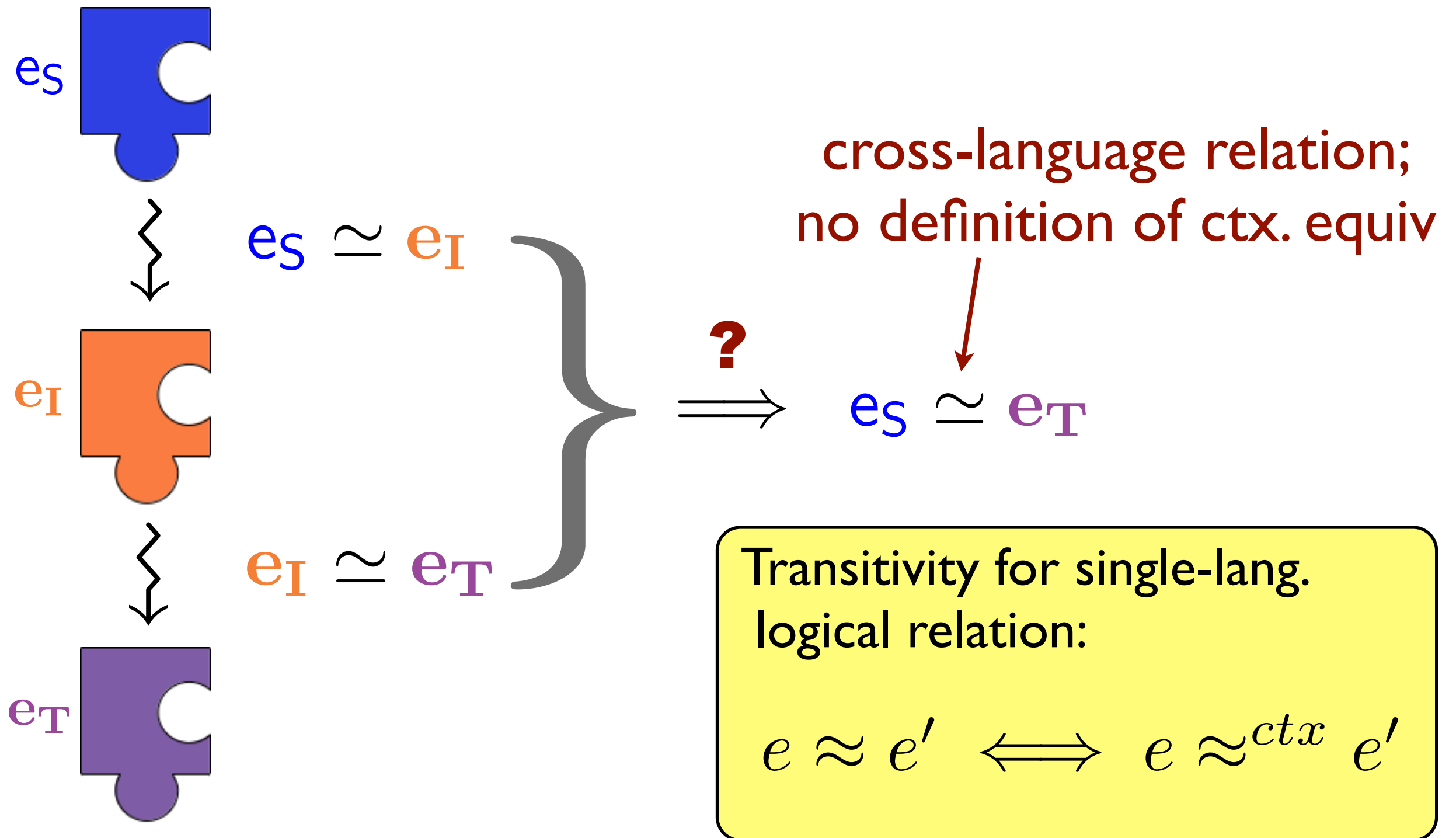
Cross-Language LR: Problem 2



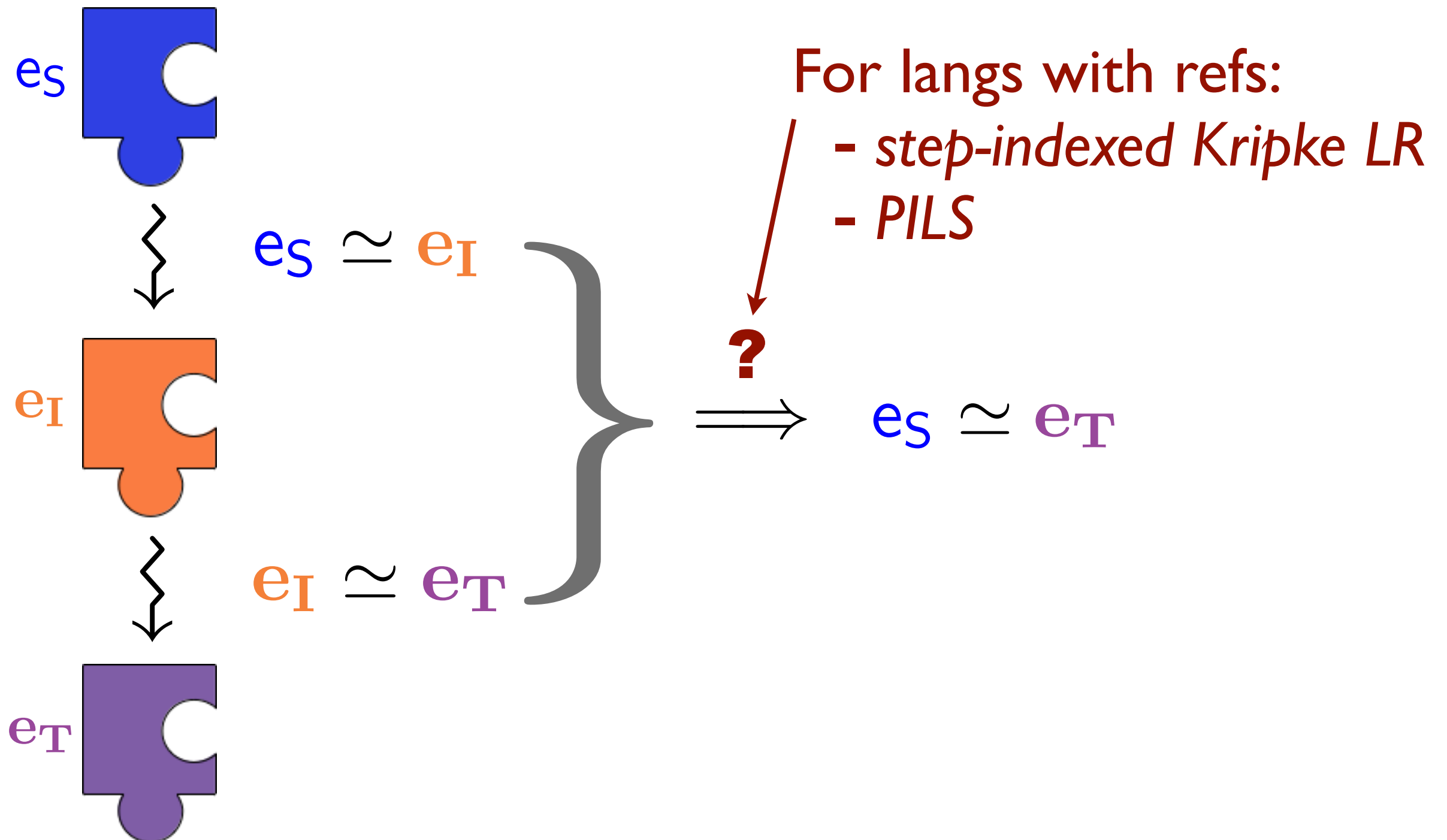
Transitivity for single-lang. logical relation?

$$\left. \begin{array}{l} e_1 \approx e_2 \\ e_2 \approx e_3 \end{array} \right\} \stackrel{?}{\implies} e_1 \approx e_3$$

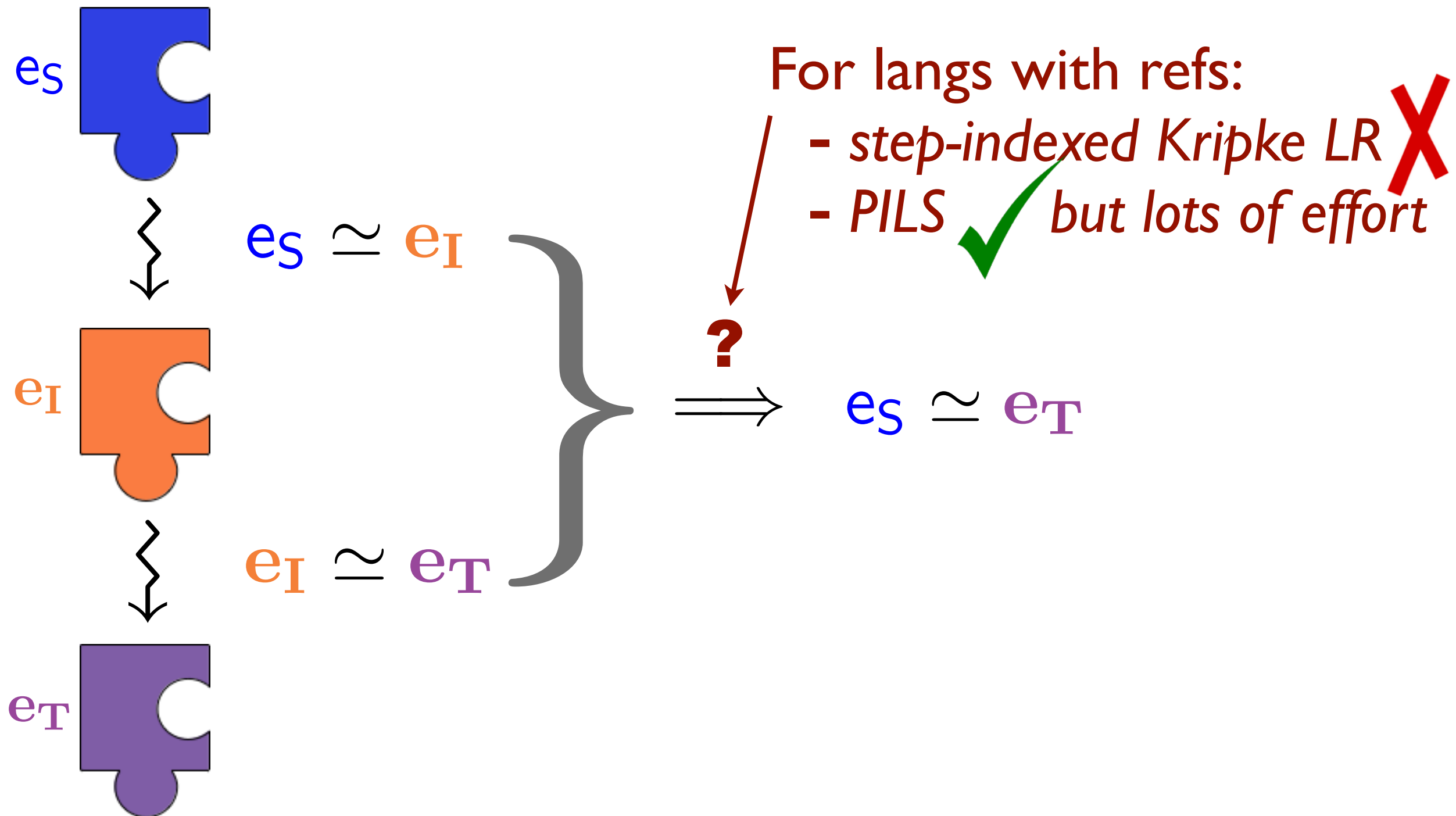
Cross-Language LR: Problem 2



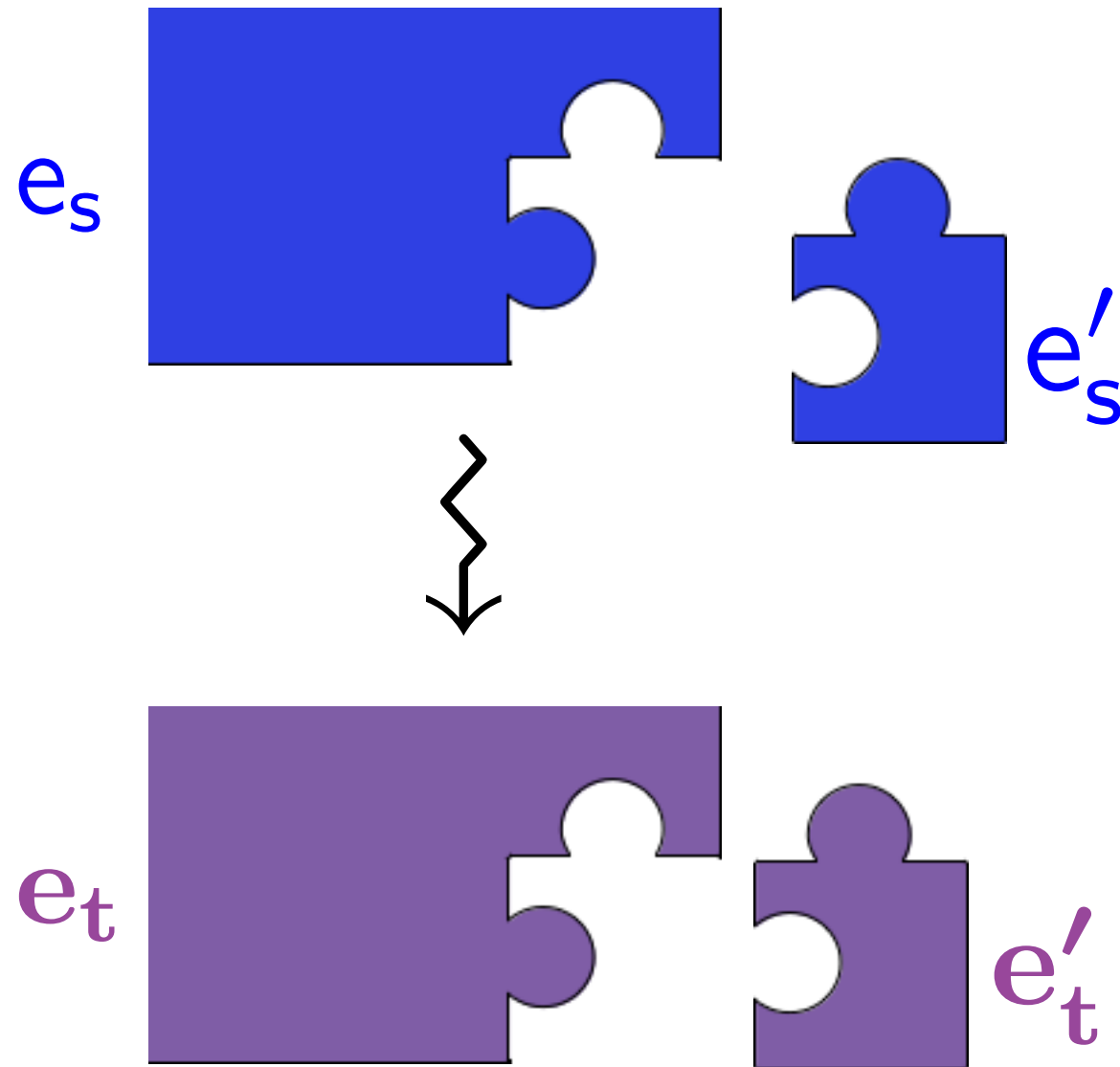
Cross-Language LR: Problem 2



Cross-Language LR: Problem 2



PILS: Problem 1 remains



Have $x : \tau' \vdash e_s \simeq e_t : \tau$

↑
PILS

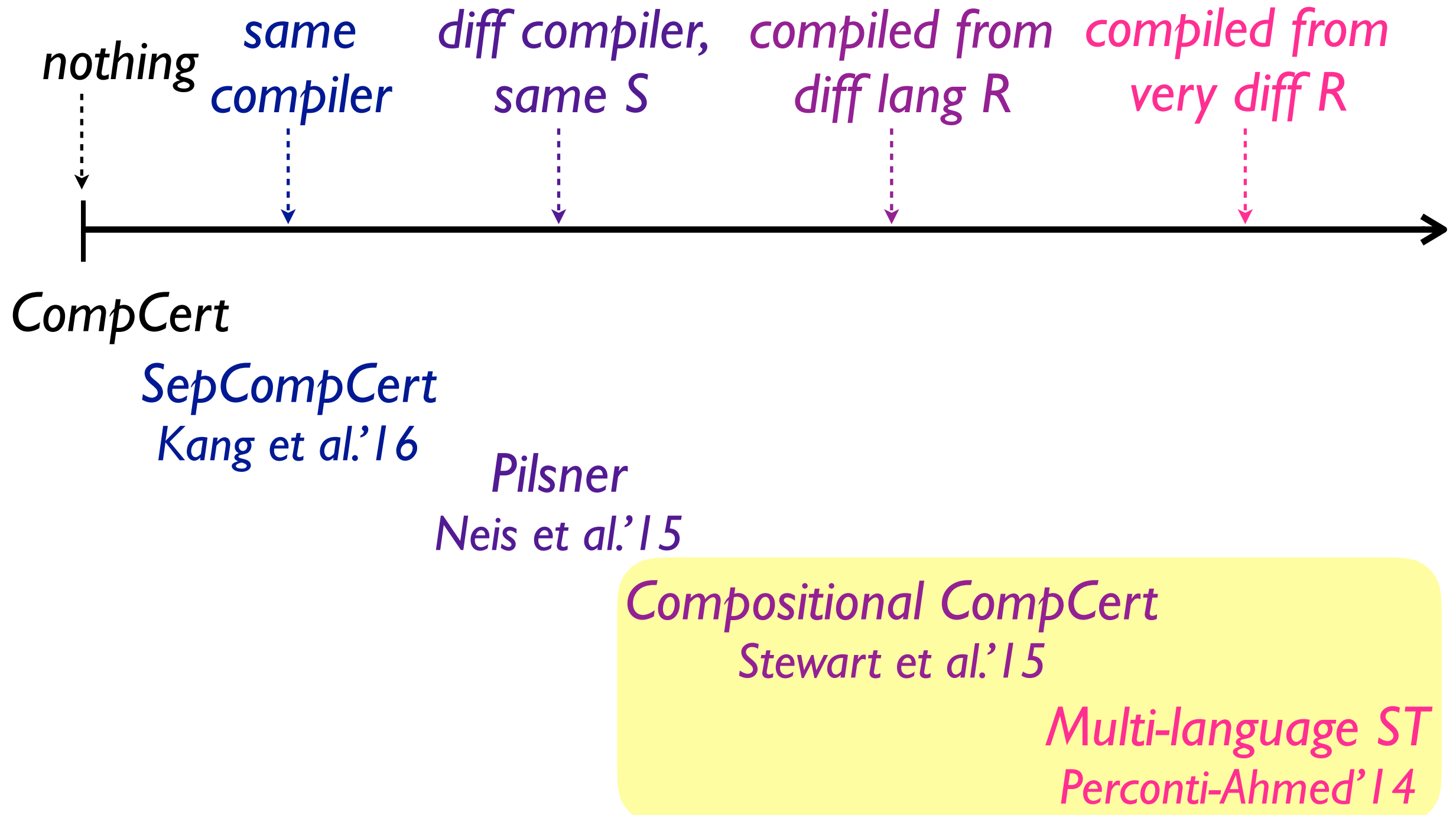
$\vdash e'_s \simeq e'_t : \tau'$

- Need to come up with e'_s
-- not feasible in practice!
- Cannot link with e'_t
whose behavior cannot
be expressed in source.

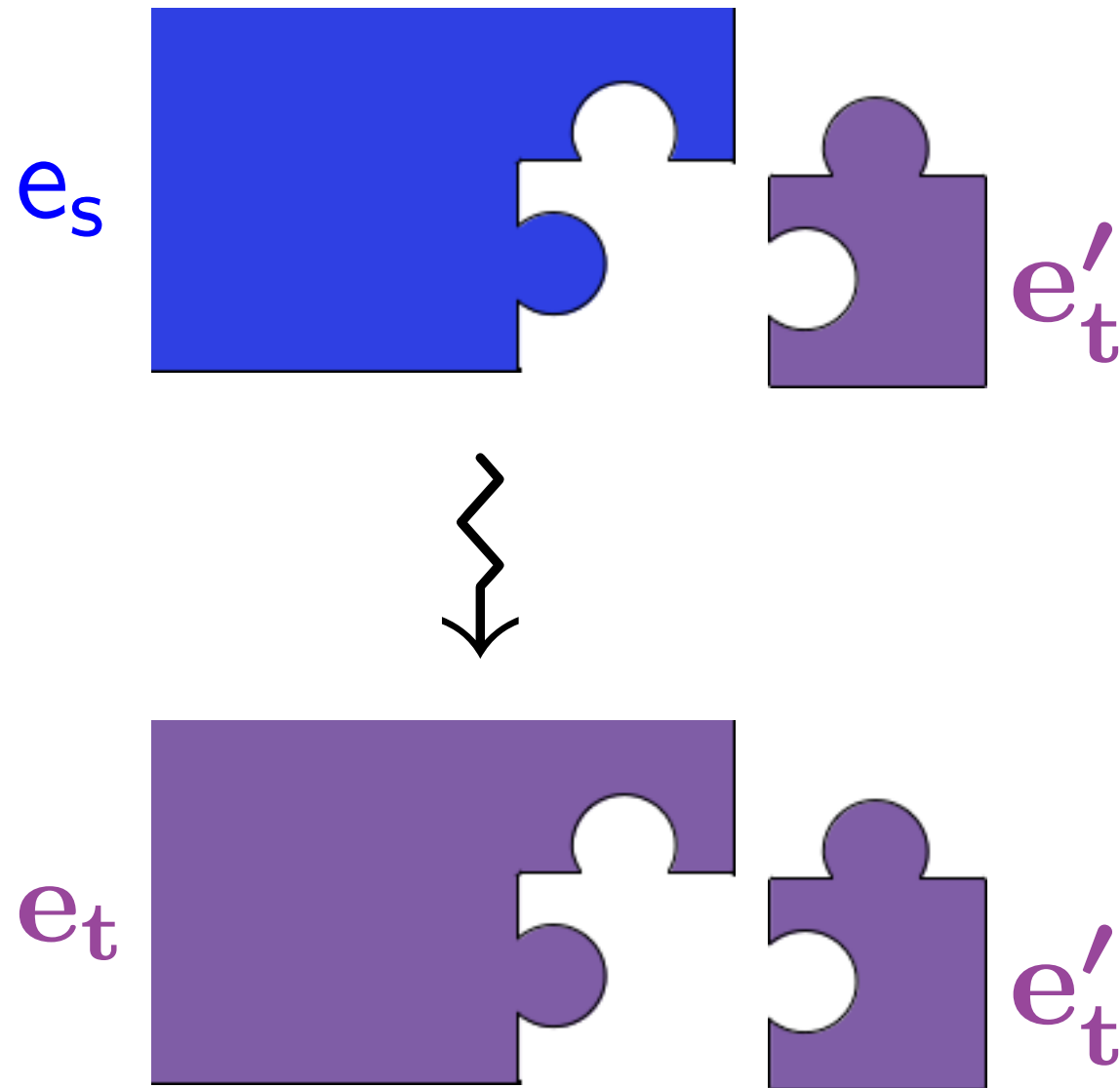
Need a New Approach...

- that works for multi-pass compilers
- that allows linking with target code of arbitrary provenance

What we can link with



Correct Compilation of Components?



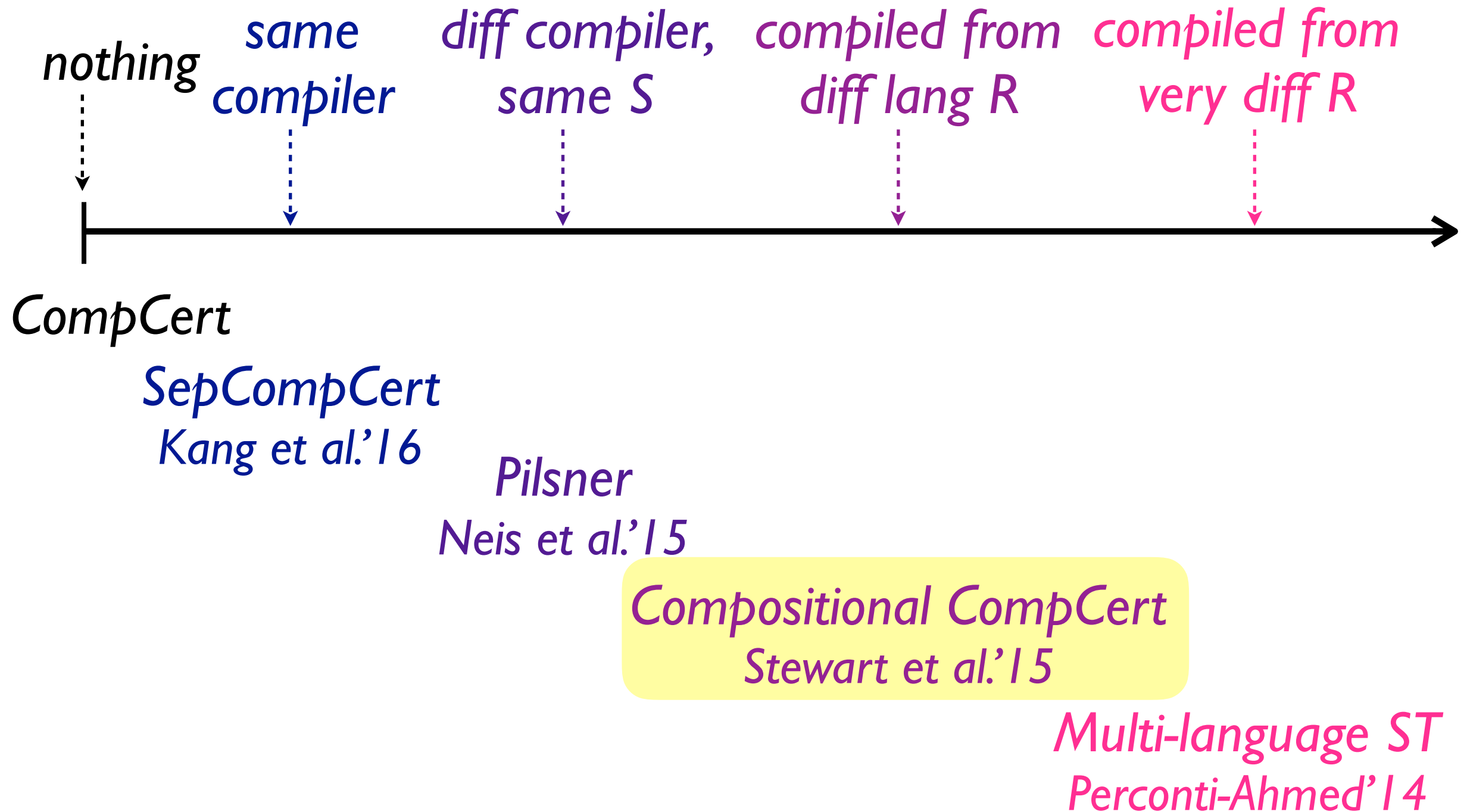
$$e_s \approx e_T$$

expressed how?

Need a semantics of source-target interoperability:

- *interaction semantics*
- *source-target multi-language*

What we can link with

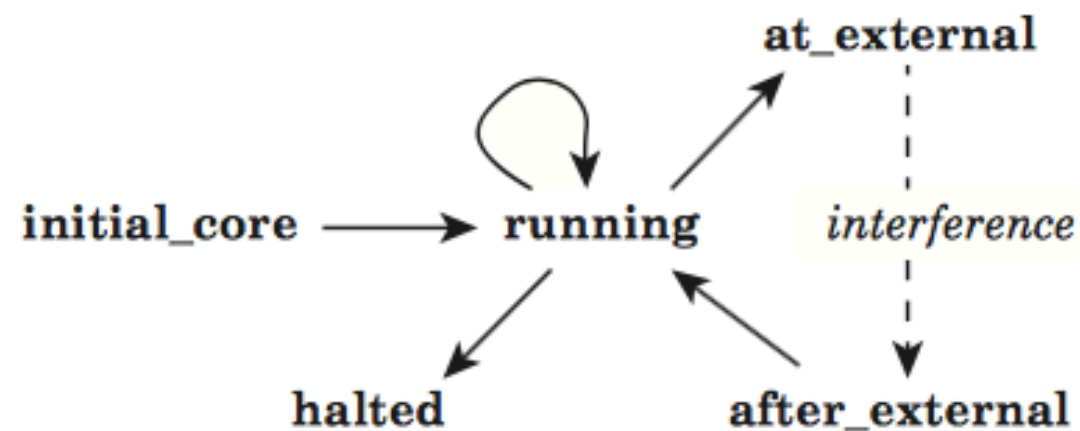


Approach: Interaction Semantics

Compositional CompCert

[Stewart et al. '15]

- Language-independent linking



```
Semantics (G C M : Type) : Type  $\triangleq$   
{  
  initial_core   : G  $\rightarrow$   $\mathcal{V}$   $\rightarrow$  list  $\mathcal{V}$   $\rightarrow$  option C  
  at_external   : C  $\rightarrow$  option ( $\mathcal{F}$   $\times$  list  $\mathcal{V}$ )  
  after_external : option  $\mathcal{V}$   $\rightarrow$  C  $\rightarrow$  option C  
  halted       : C  $\rightarrow$  option  $\mathcal{V}$   
  corestep     : G  $\rightarrow$  C  $\rightarrow$  M  $\rightarrow$  C  $\rightarrow$  M  $\rightarrow$  Prop  
}
```

Figure 2. Interaction semantics interface. The types G (global environment), C (core state), and M (memory) are parameters to the interface. \mathcal{F} is the type of external function identifiers. \mathcal{V} is the type of CompCert values.

Approach: Interaction Semantics

Compositional CompCert

[Stewart et al. '15]

- Language-independent linking
- Structured simulation: support rely-guarantee relationship between the different languages while retaining vertical compositionality

Semantic representation of contexts code can link with.

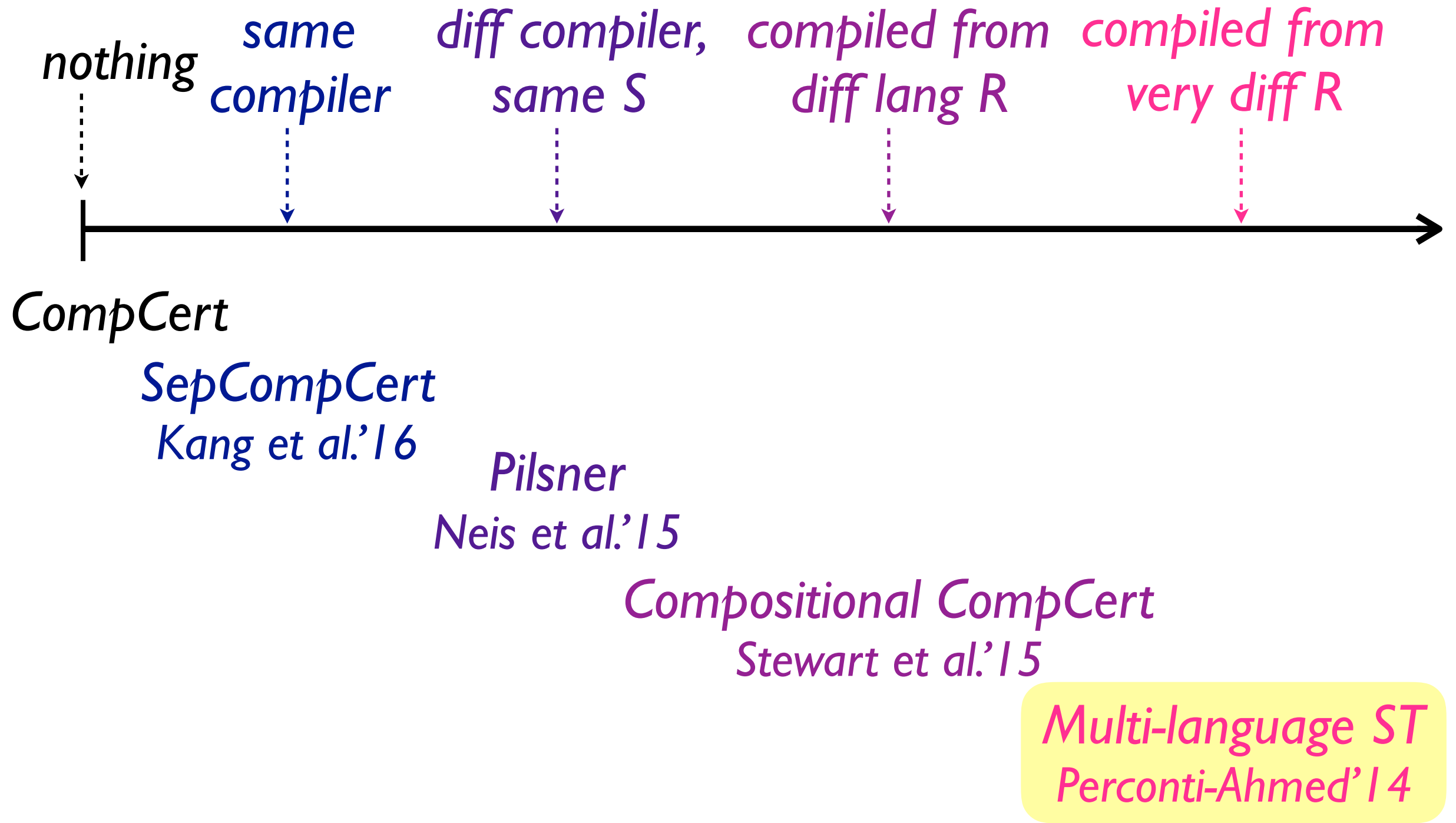
Approach: Interaction Semantics

Compositional CompCert

[Stewart et al. '15]

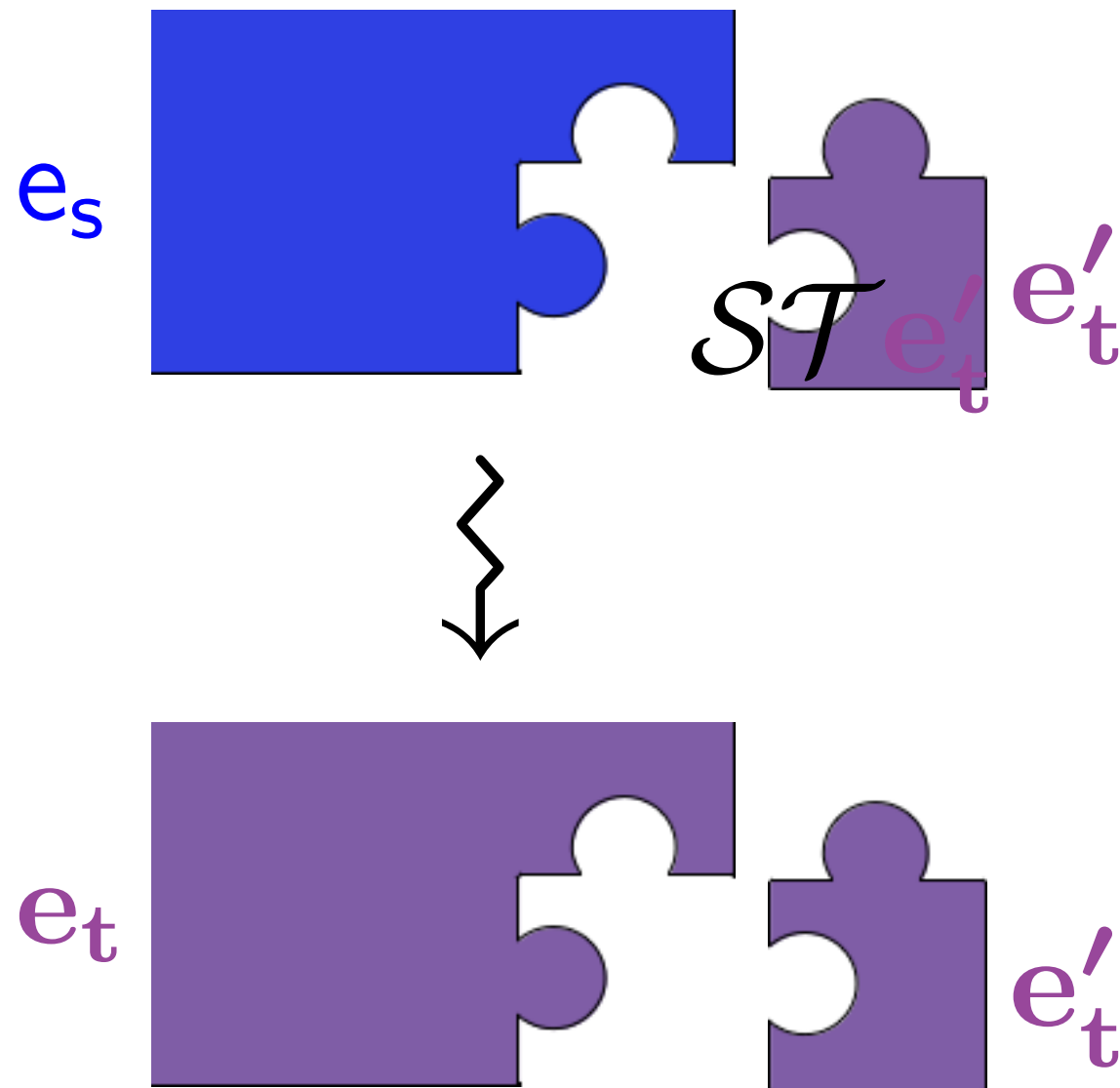
- Language-independent linking
 - uniform CompCert memory model across all languages
 - not clear how to scale to richer source langs (e.g., ML), compilers with different source/target memory models
- Structured simulation: support rely-guarantee relationship between the different languages while retaining vertical compositionality
 - transitivity relies on compiler passes performing restricted set of memory transformations

What we can link with



Approach: Source-Target Multi-lang.

[Perconti-Ahmed'14]



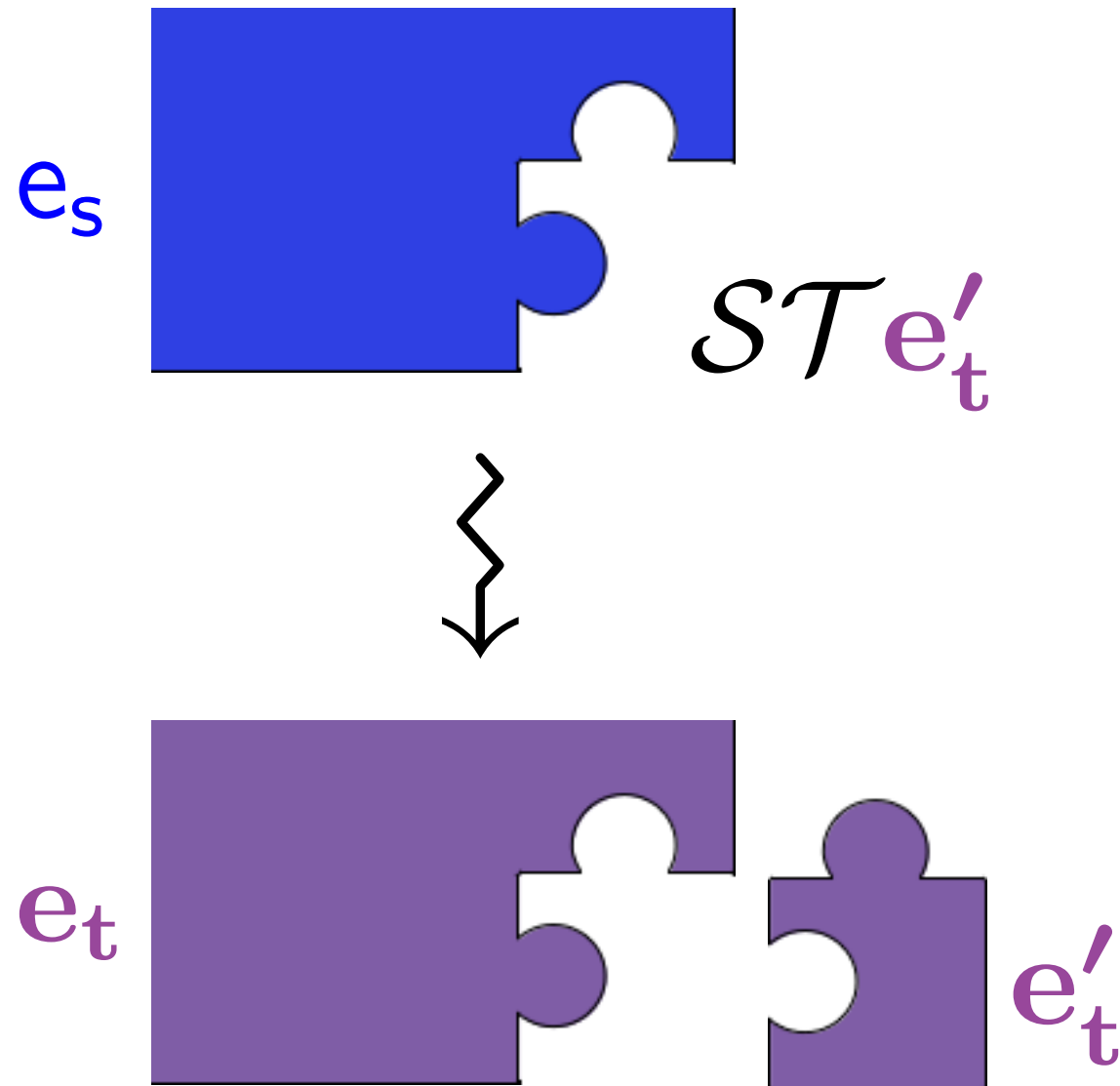
Specify semantics
of source-target
interoperability:

$ST e_t$ $TS e_s$

*Multi-language semantics:
a la Matthews-Findler '07*

Approach: Source-Target Multi-lang.

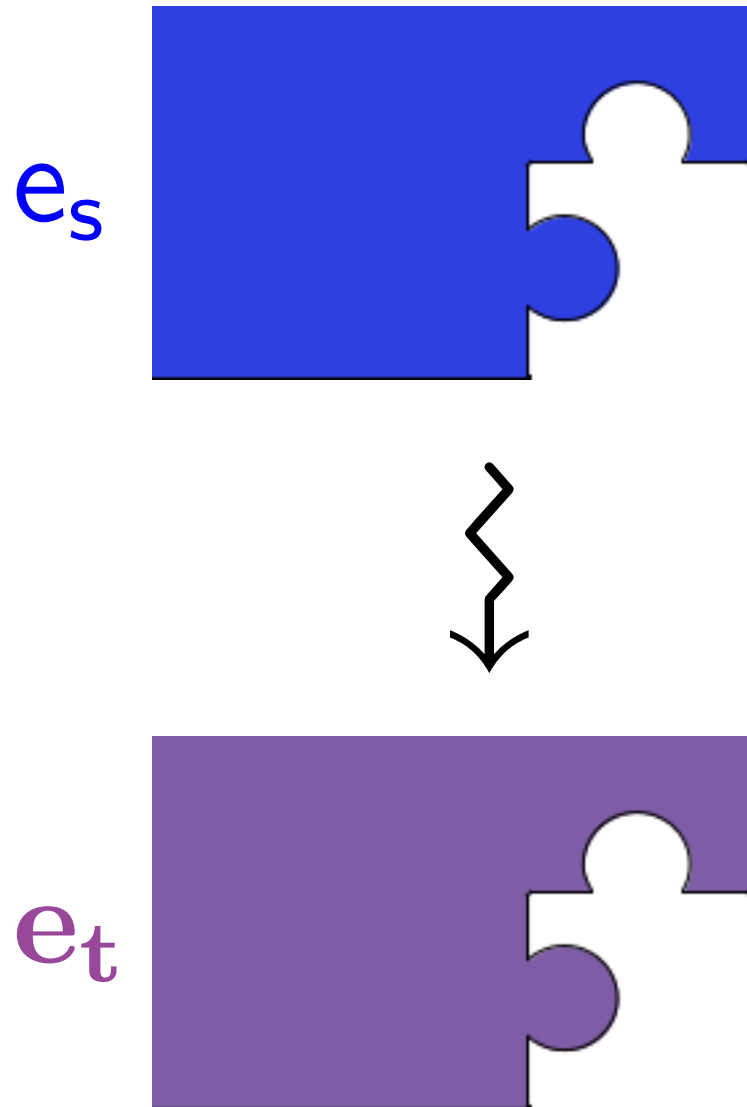
[Perconti-Ahmed'14]



$$\mathcal{TS}(e_s (ST e'_t)) \approx^{ctx} e_t e'_t$$

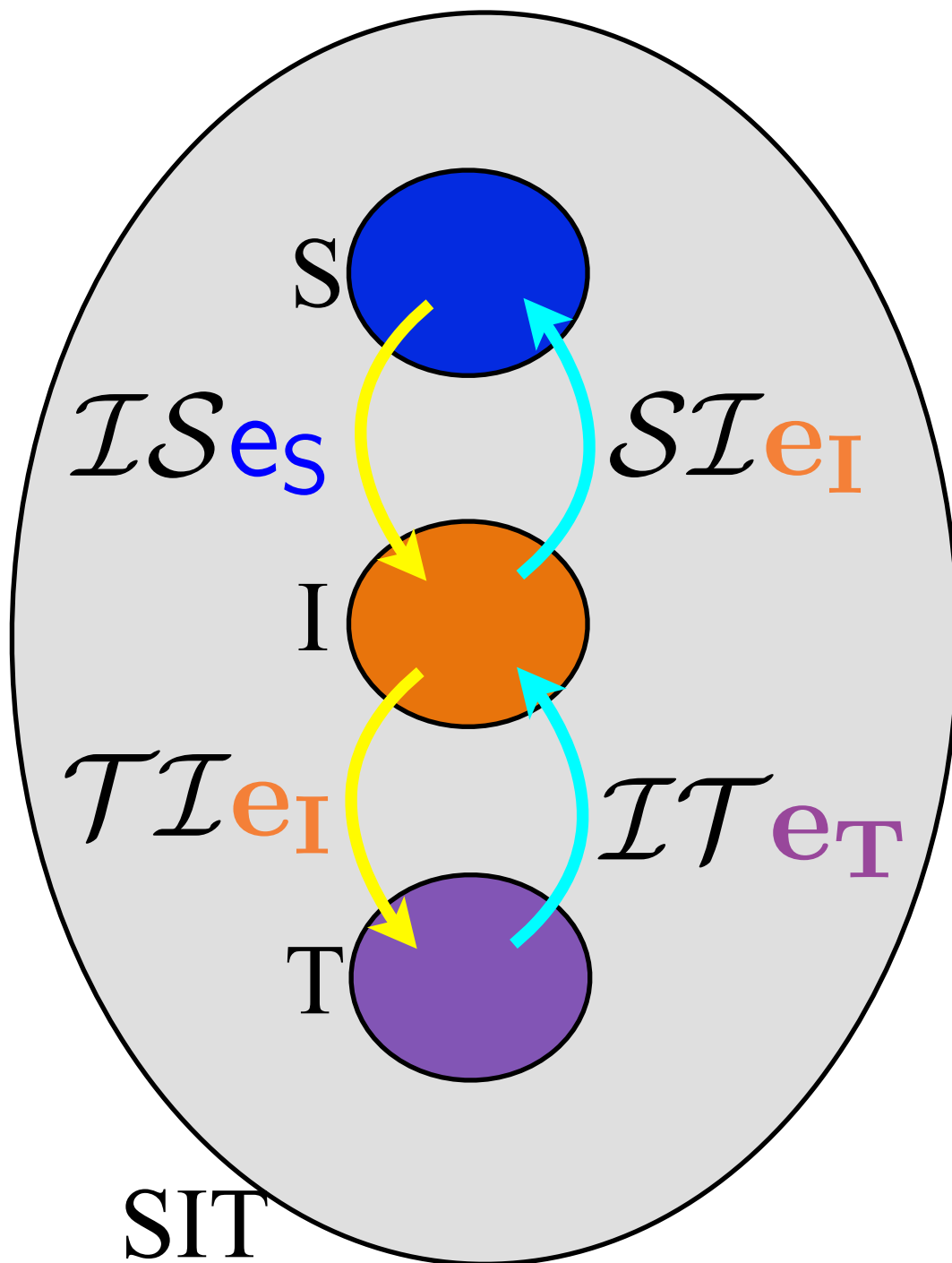
Approach: Source-Target Multi-lang.

[Perconti-Ahmed'14]

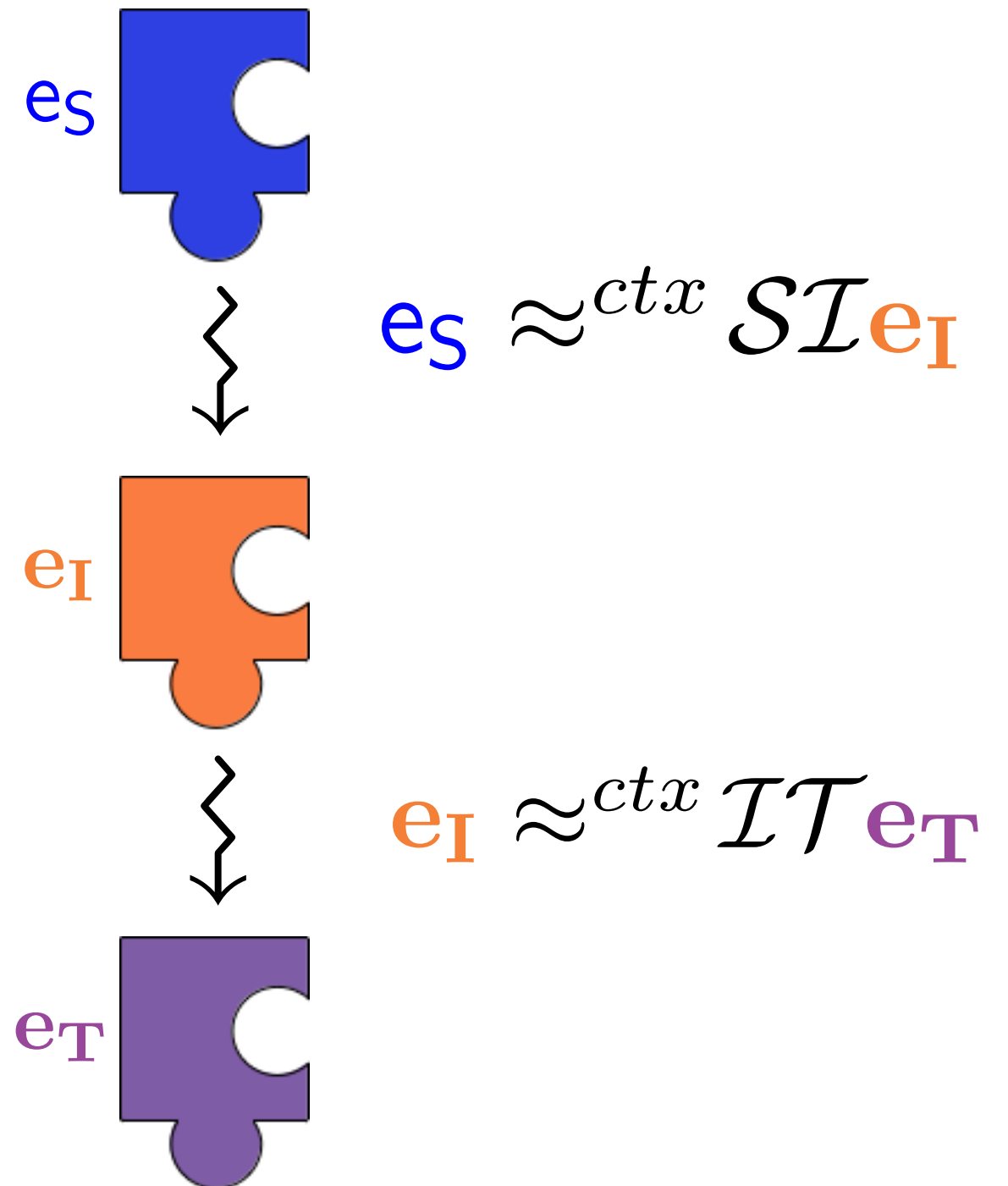


$$e_s \approx e_T \stackrel{\text{def}}{=} e_s \approx^{ctx} \mathcal{ST} e_T$$

Multi-Language Semantics Approach

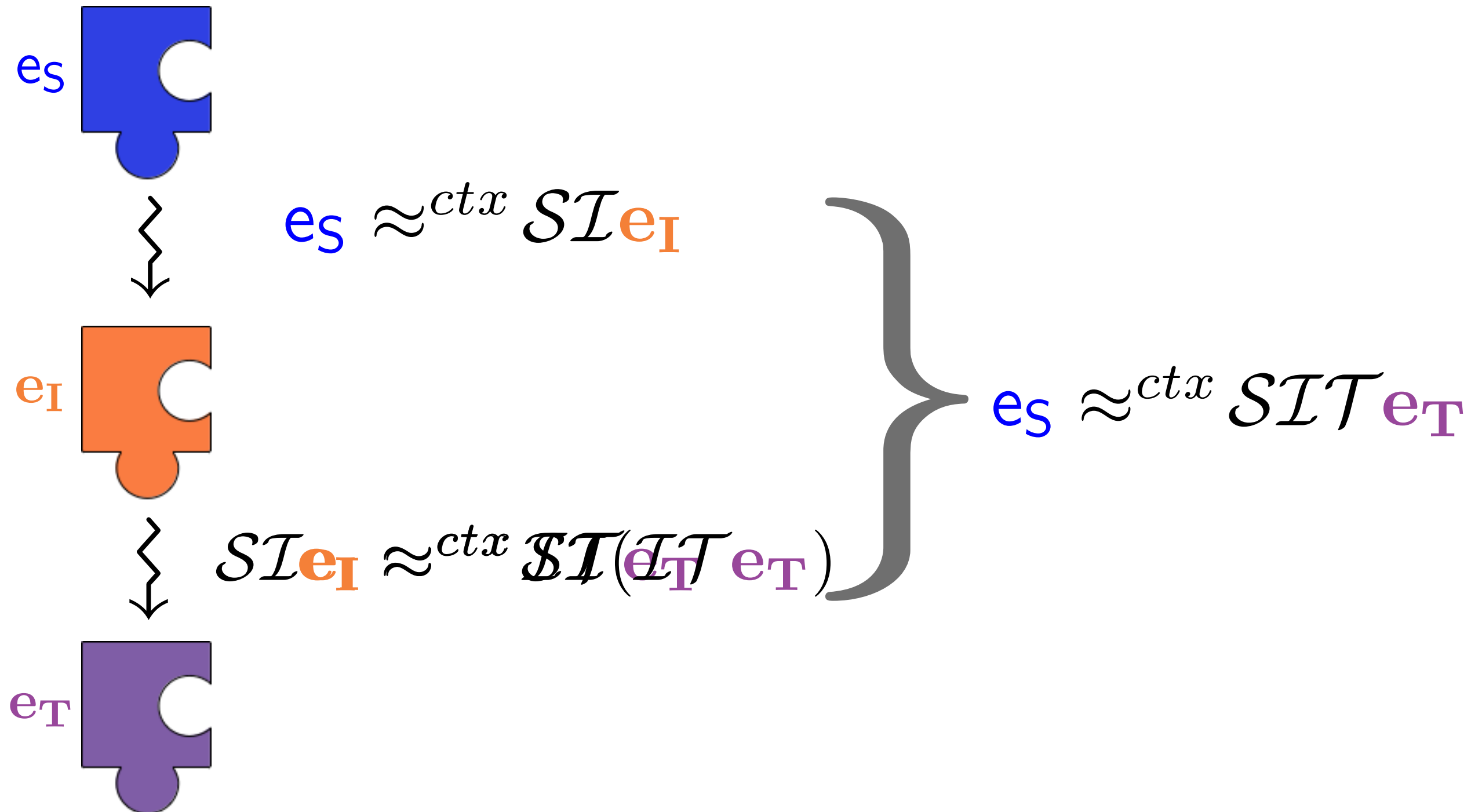


Compiler Correctness

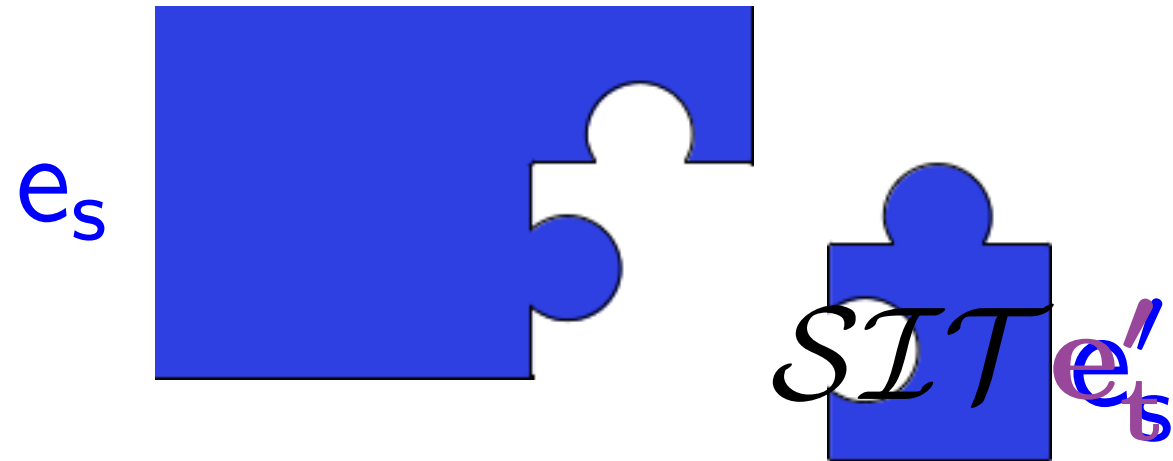


Multi-Lang. Approach: Multi-pass ✓

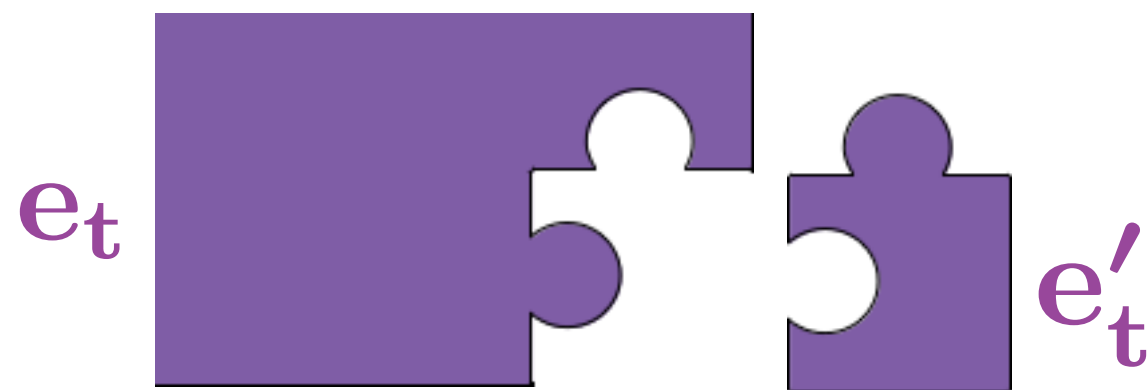
Compiler Correctness



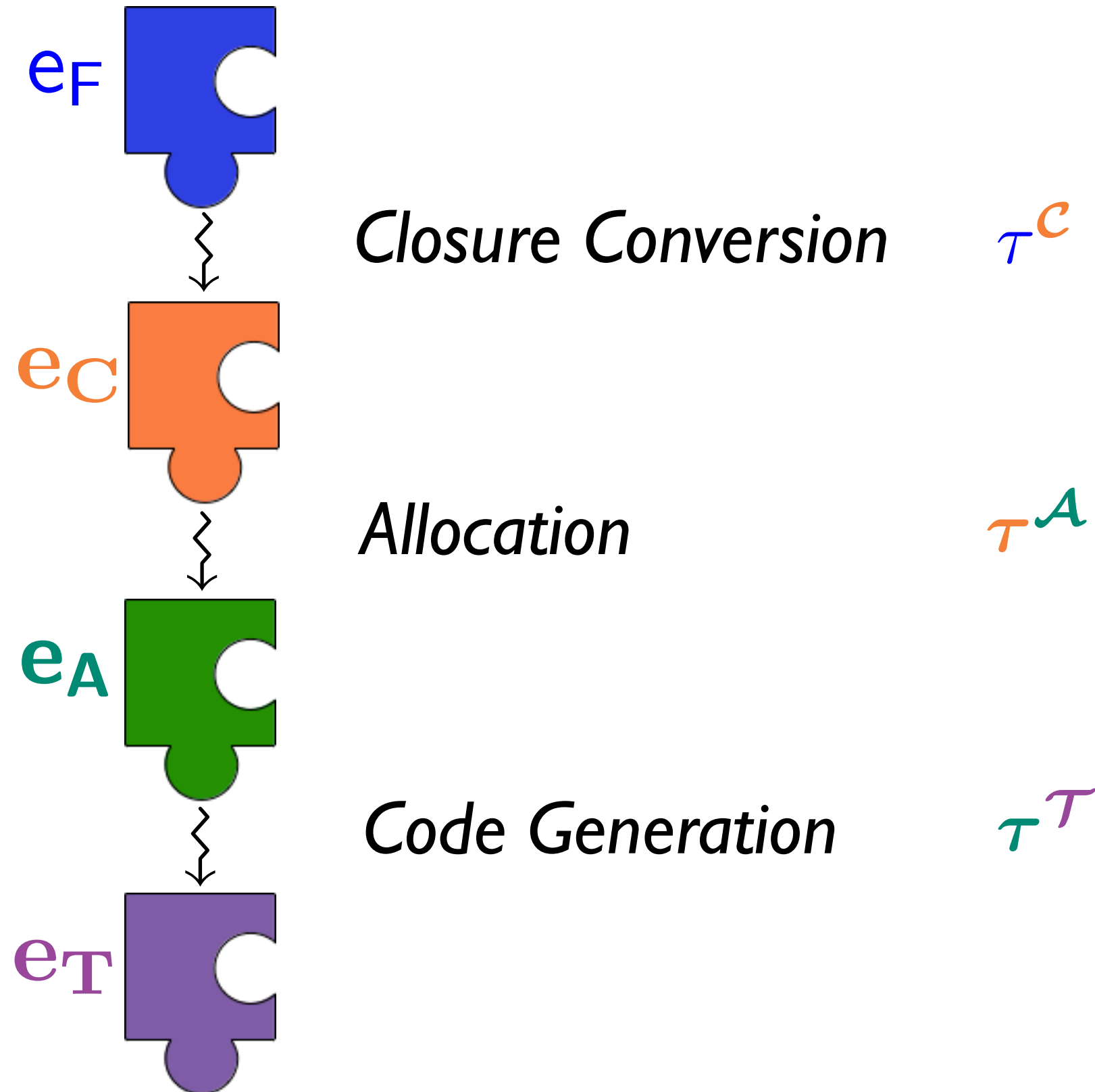
Multi-Lang. Approach: Linking ✓



$$TIS(e_s (SIT e'_s)) \approx^{ctx} e_t e'_t$$



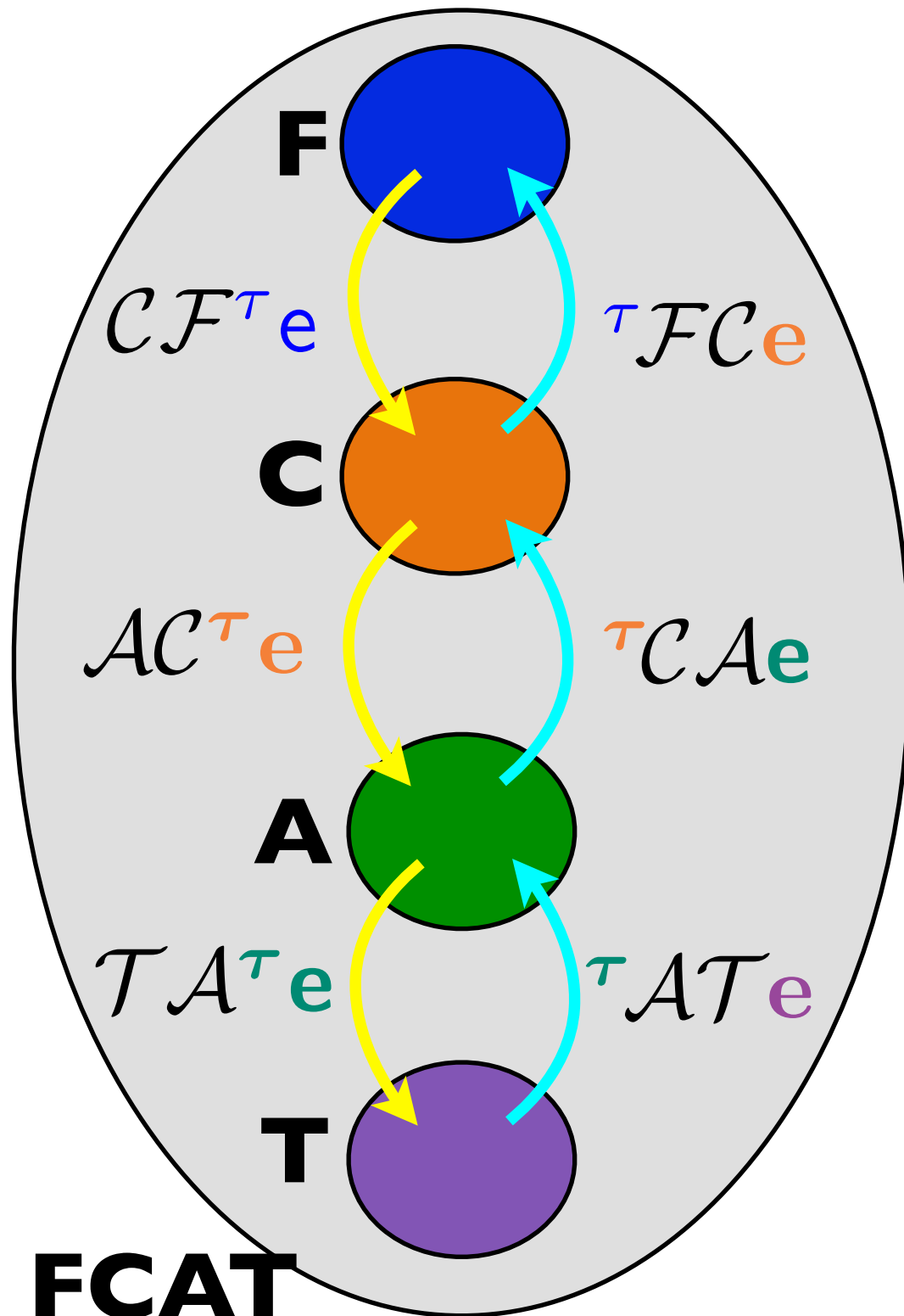
Compiler Correctness: F to TAL



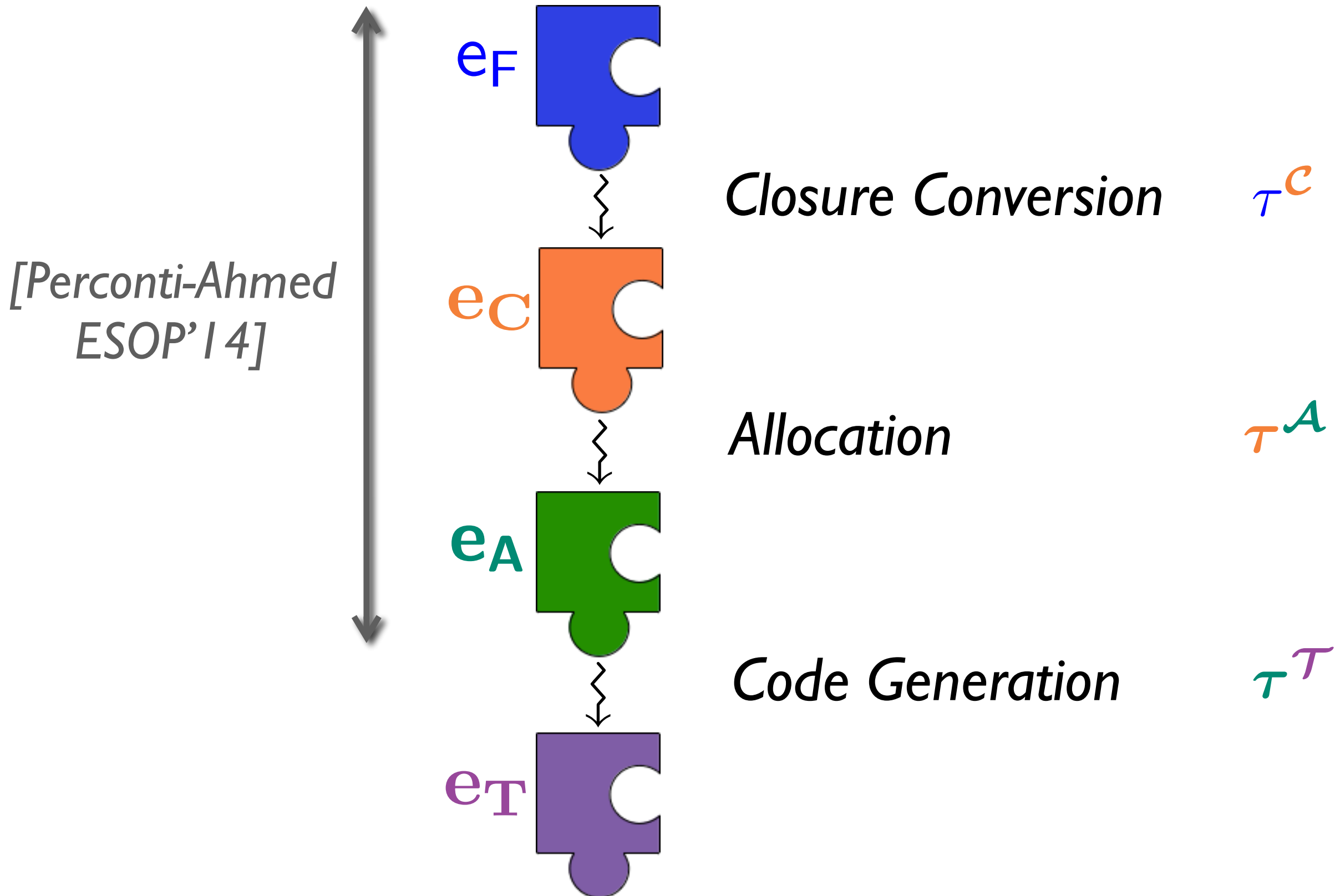
Combined language **FCAT**

- Boundaries mediate between

τ & τ^C τ & τ^A τ & τ^T



Compiler Correctness: F to TAL



CompCompCert vs. Multi-language

Transitivity:

- structured simulations

- all passes use multi-lang \approx^{ctx}

Check okay-to-link-with:

- satisfies CompCert
memory model

- satisfies expected type
(translation of source type)

Contexts:

- semantic representation

- syntactic representation

Requires uniform memory model across compiler IRs?

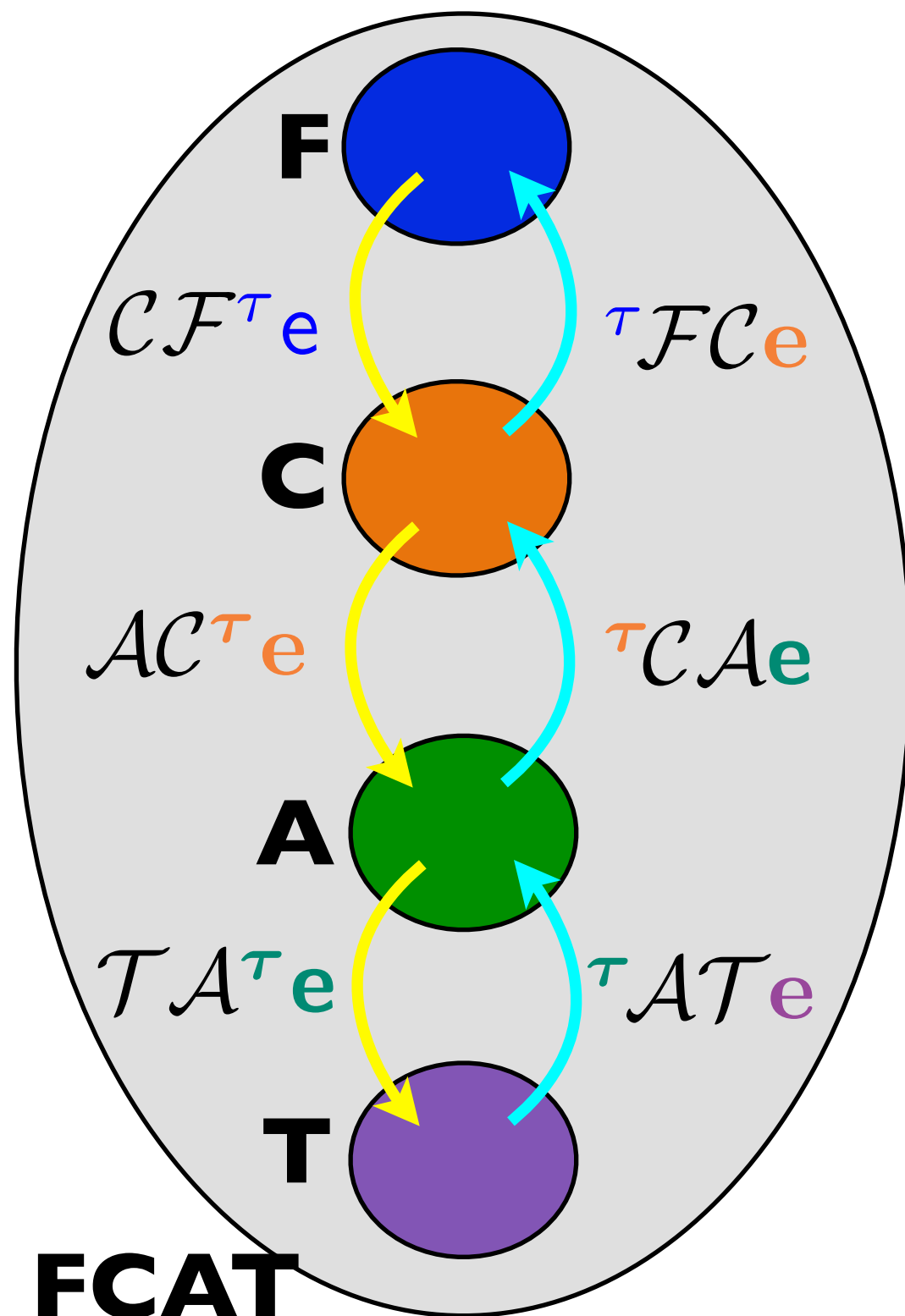
- yes

- no

Case Study: Closure Conversion

Correctness of typed closure conversion
using multi-language semantics... [on board]

Challenges

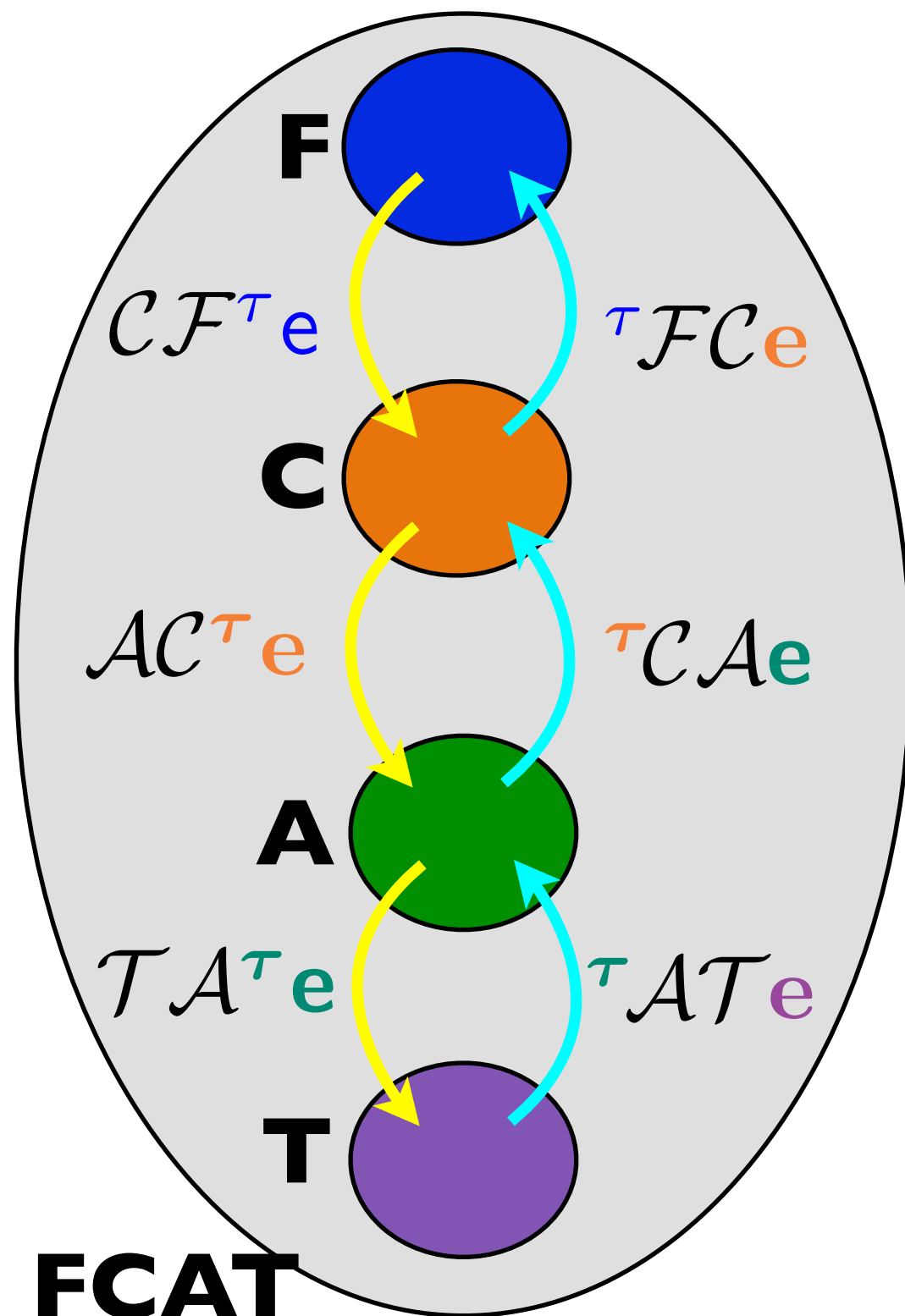


F+C: Interoperability semantics with type abstraction in both languages

C+A: Interoperability when compiler pass allocates code & tuples on heap

A+T: What is e ? What is v ?
How to define contextual equiv. for TAL *components*?
How to define logical relation?

Challenges



F+C: Interoperability semantics with type abstraction in both languages

C+A: Interoperability when compiler pass allocates code & tuples on heap

A+T: What is e ? What is v ?
How to define contextual equiv. for TAL *components*?
How to define logical relation?

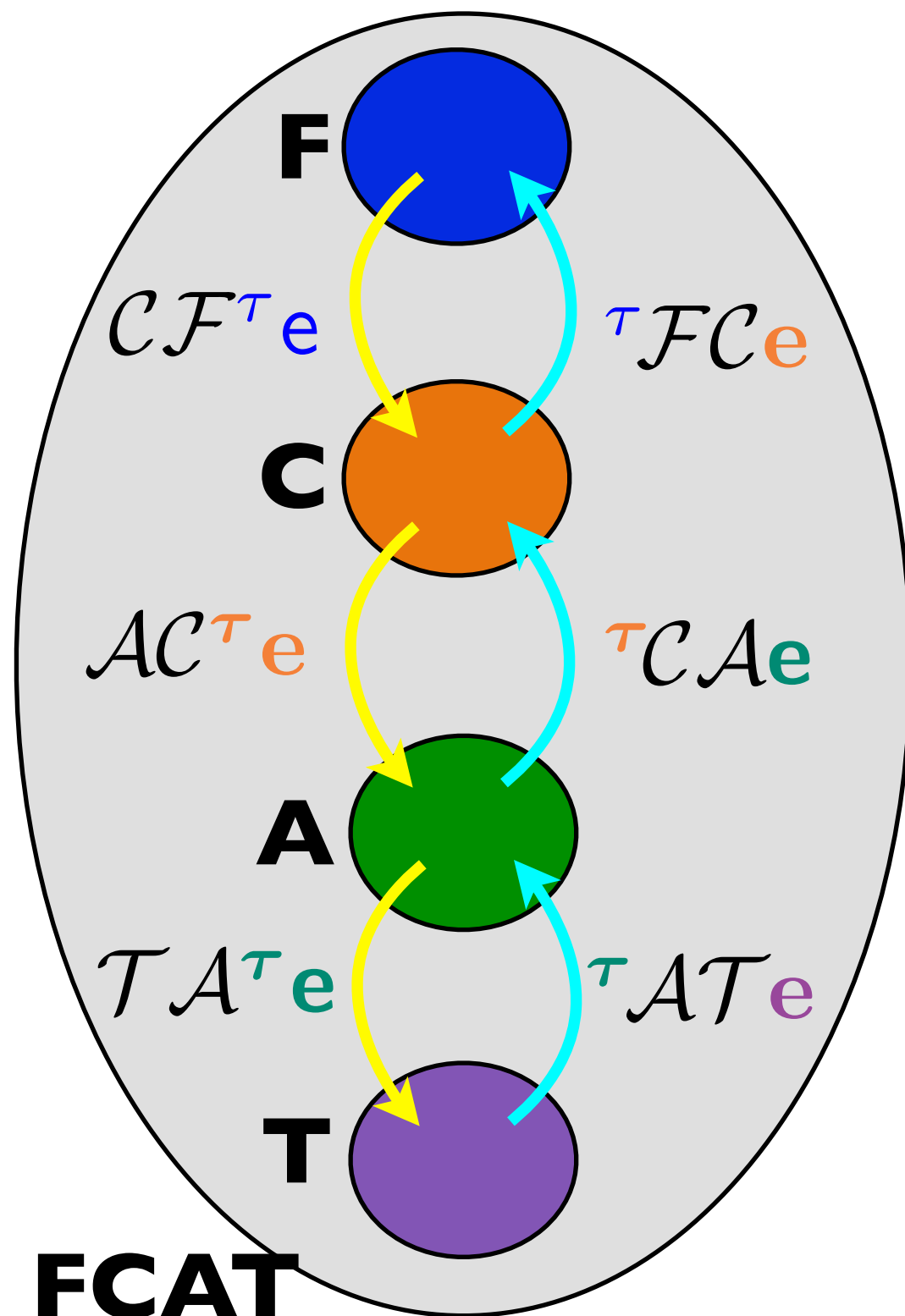
Interoperability: **C** and **A**

$$\mathbf{H}, \ell \mapsto \langle \mathbf{v} \rangle; \langle \tau \rangle \mathbf{C} \mathbf{A} \ell \mapsto \mathbf{H}, \ell \mapsto \langle \mathbf{v} \rangle; \langle \mathbf{v} \rangle$$

$$\mathbf{H}; \mathbf{A} \mathbf{C} \langle \tau \rangle \langle \mathbf{v} \rangle \mapsto \mathbf{H}, \ell \mapsto \langle \mathbf{v} \rangle; \ell$$

Allocate a new
location for tuple

Challenges



F+C: Interoperability semantics with type abstraction in both languages

C+A: Interoperability when compiler pass allocates code & tuples on heap

A+T: What is e ? What is v ?
How to define contextual equiv. for TAL components?
How to define logical relation?

Central Challenge: interoperability between
high-level (direct-style) language &
assembly (continuation style)

FunTAL: Reasonably Mixing a Functional Language
with Assembly [*Patterson et al. '17*]

What is a component in TAL?

$$e : \mathcal{T} \rightsquigarrow e$$

e



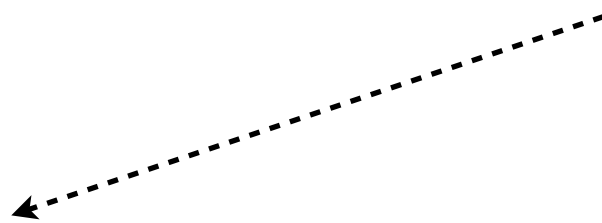
Instruction Sequence



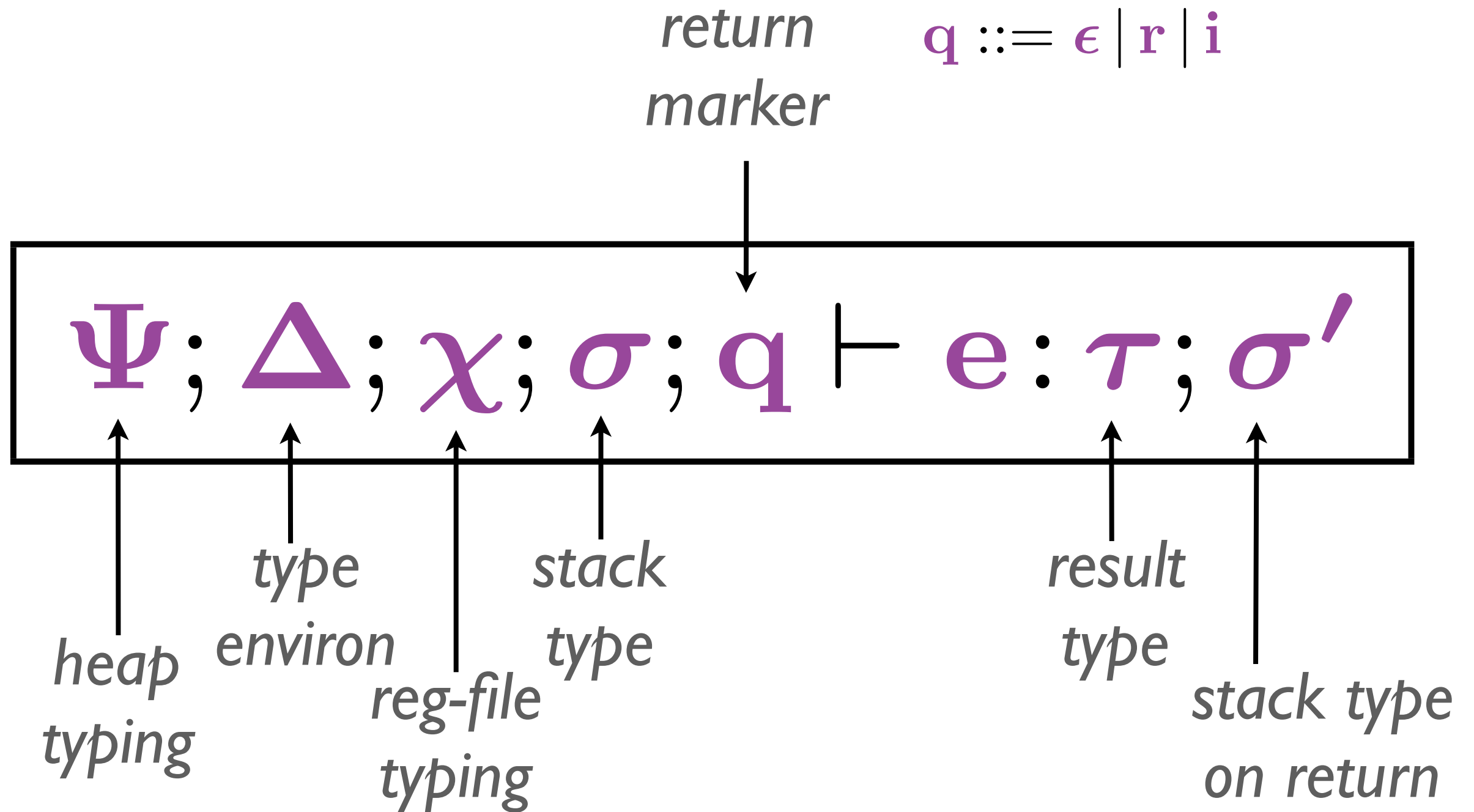
$$e ::= (I, H)$$



Heap with basic blocks



Typing TAL Components



Basic blocks

e



$$e ::= (\mathbf{I}, \mathbf{H})$$


$$\text{code}[\Delta]\{\chi; \sigma\}^q.\mathbf{I} : \forall[\Delta].\{\chi; \sigma\}^q$$

Equivalence of \mathbf{T} Components: Tricky!

Logical relations: related inputs to related outputs

$$\mathcal{V}[\tau_1 \rightarrow \tau_2] = \{(W, \lambda x.e_1, \lambda x.e_1) \mid \dots\}$$

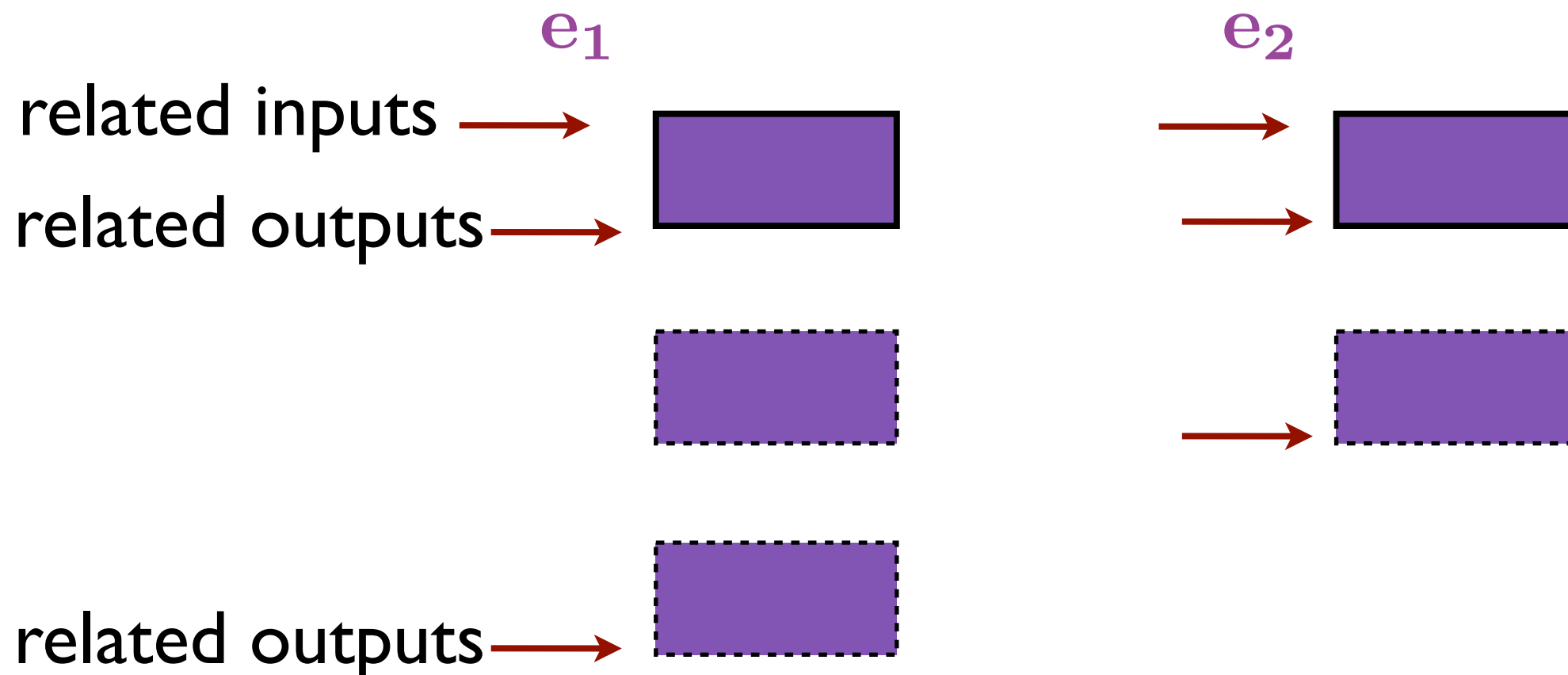
$$\mathcal{HV}[\forall[\Delta].\{\chi; \sigma\}^q] = \{(W, \text{code}[\Delta]\{\chi; \sigma\}^q.I_1, \text{code}[\Delta]\{\chi; \sigma\}^q.I_2) \mid \dots\}$$

Equivalence of \mathbf{T} Components

Logical relations: related inputs to related outputs

$$\mathcal{V}[\tau_1 \rightarrow \tau_2] = \{(W, \lambda x.e_1, \lambda x.e_1) \mid \dots\}$$

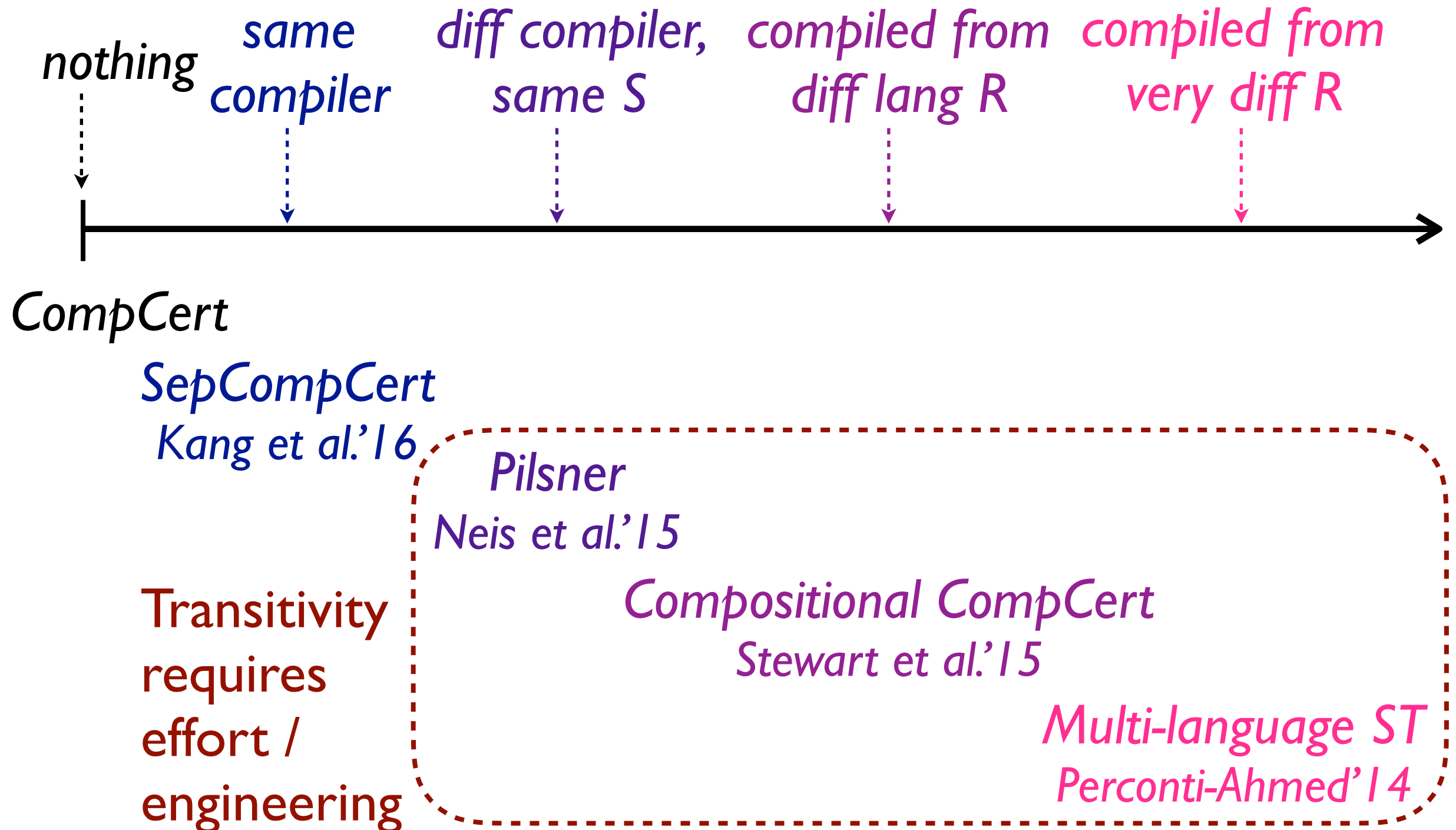
$$\mathcal{HV}[\forall[\Delta].\{\chi; \sigma\}^q] = \{(W, \text{code}[\Delta]\{\chi; \sigma\}^q.I_1, \text{code}[\Delta]\{\chi; \sigma\}^q.I_2) \mid \dots\}$$



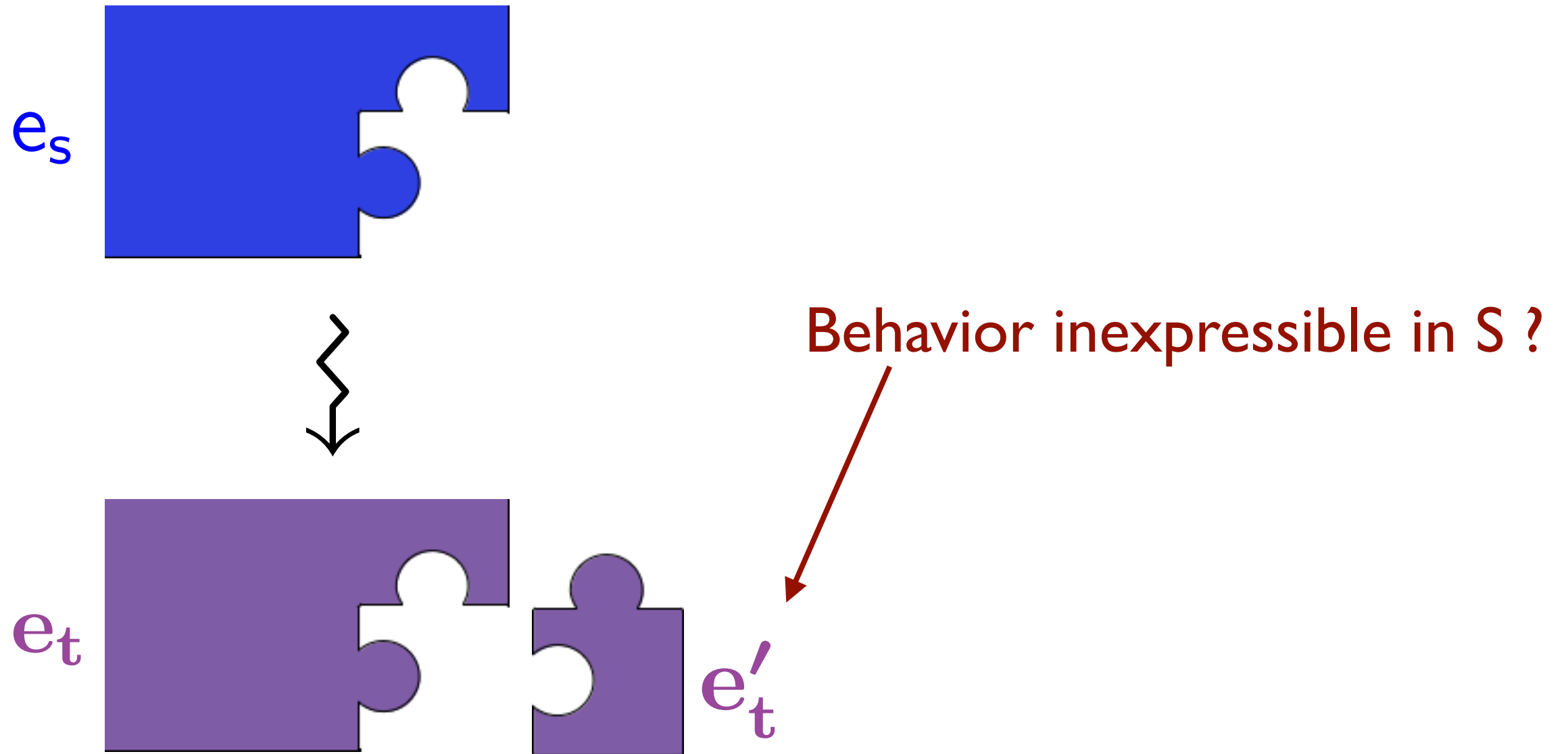
Ongoing: Multi-lang. Approach

- Underway: Code Generation pass to TAL
- Working on simplifying multi-language design to support easier proofs when multiple embedded languages have polymorphism & refs
 - *Matthew Kolosick, Dustin Jamner, Max New, AA*

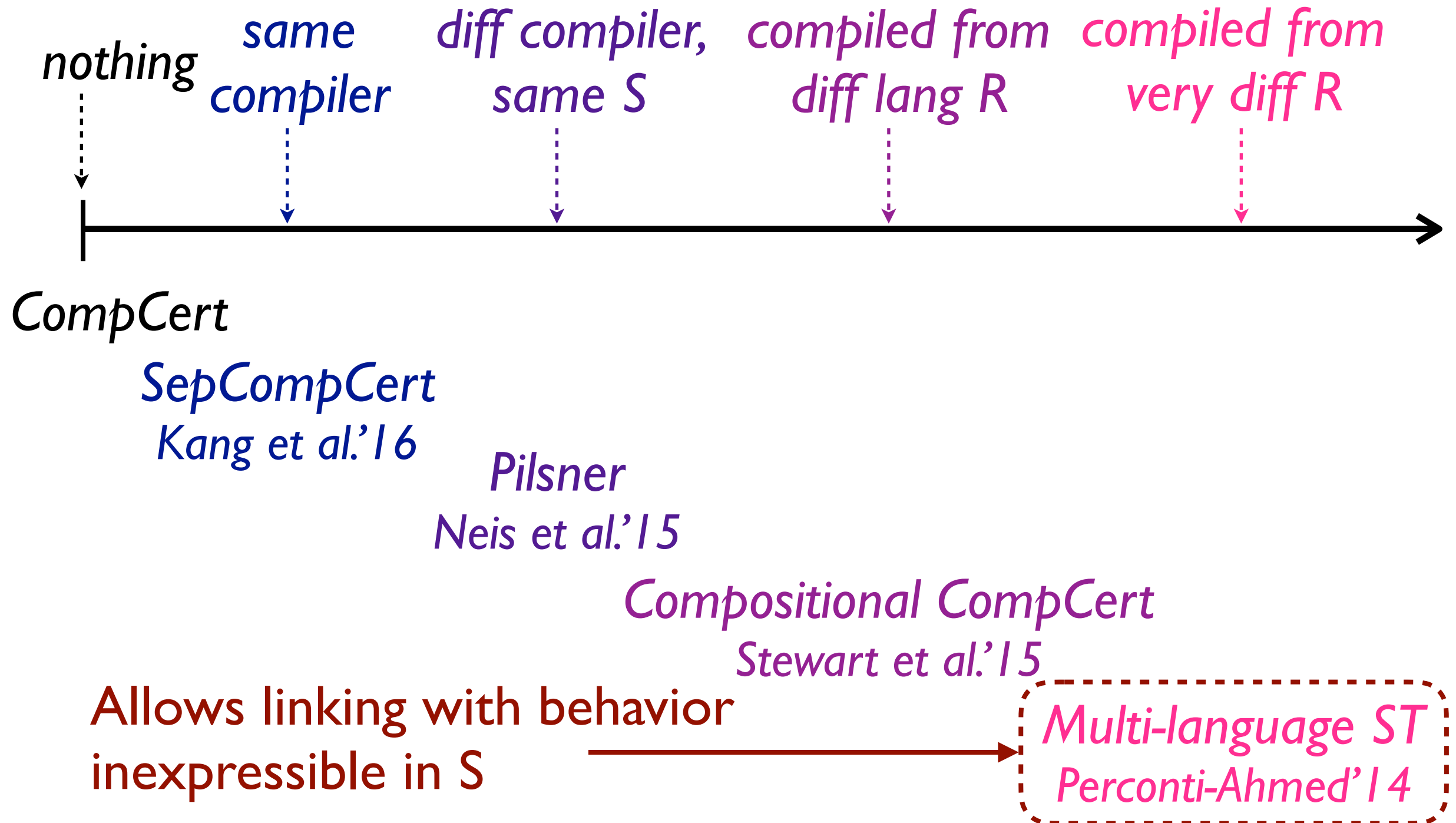
Horizontal / **Vertical** Compositionality



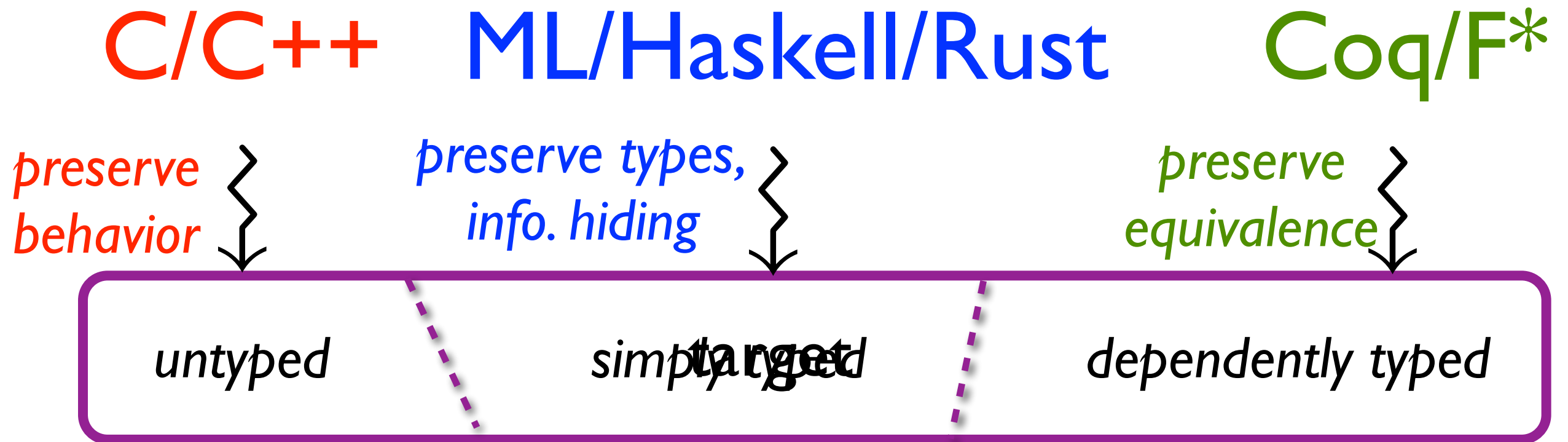
Horizontal Compositionality



Horizontal / Vertical Compositionality

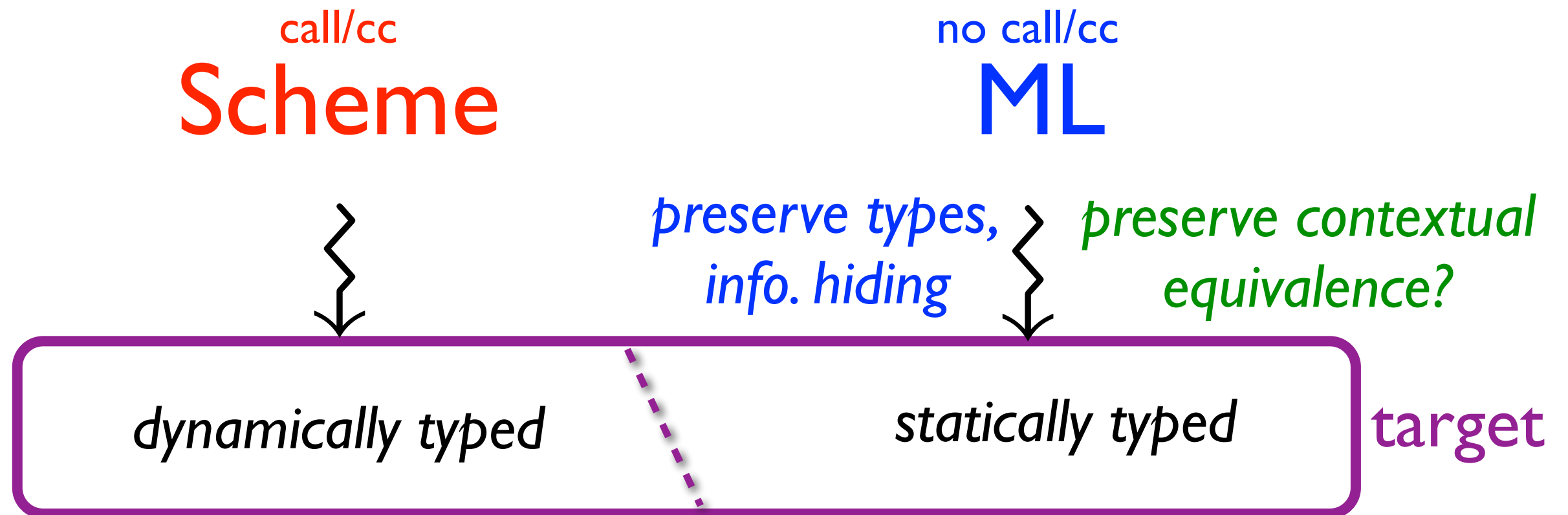


Verified Compilers for Multi-lang. World



- It's about principled language interoperability!

“Principled” Language Interoperability?



- Compiler can preserve different properties through choice of type-translation: a spectrum of linking options

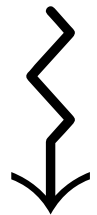
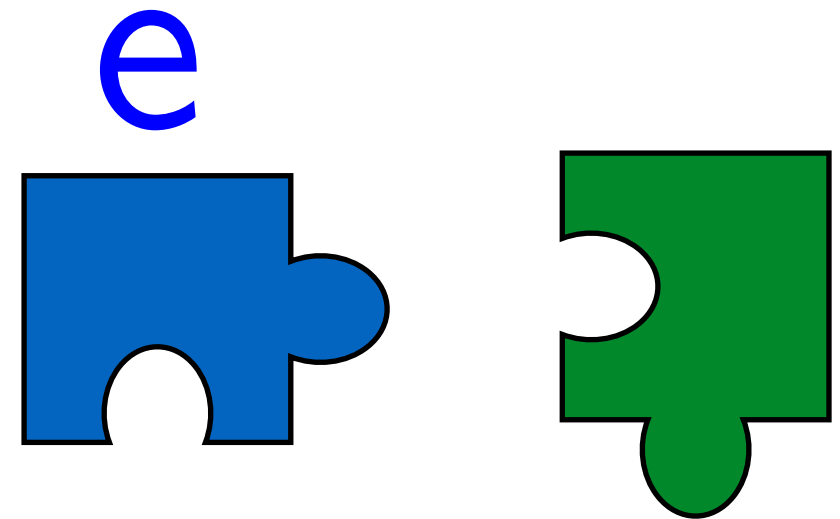
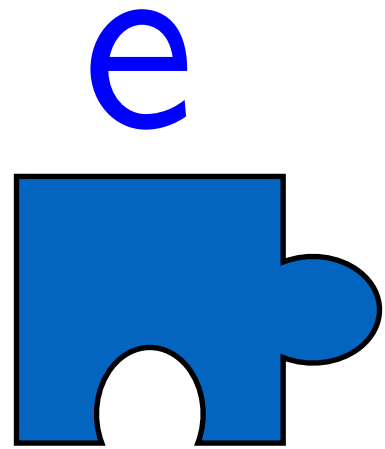
Next...

- Fully Abstract Compilation / Secure Compilation
 - compiler is equivalence-preserving
 - ensures a compiled component does not interact with any target behavior that is inexpressible in S
 - Recent results on fully abstract compilation

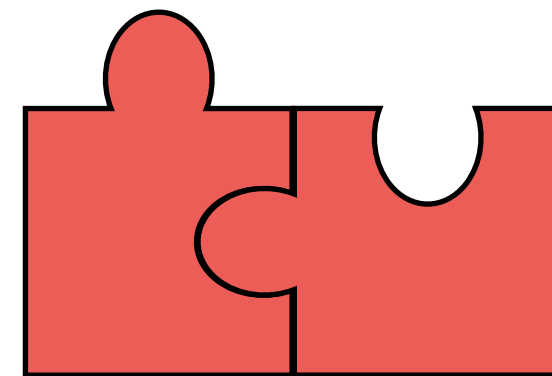
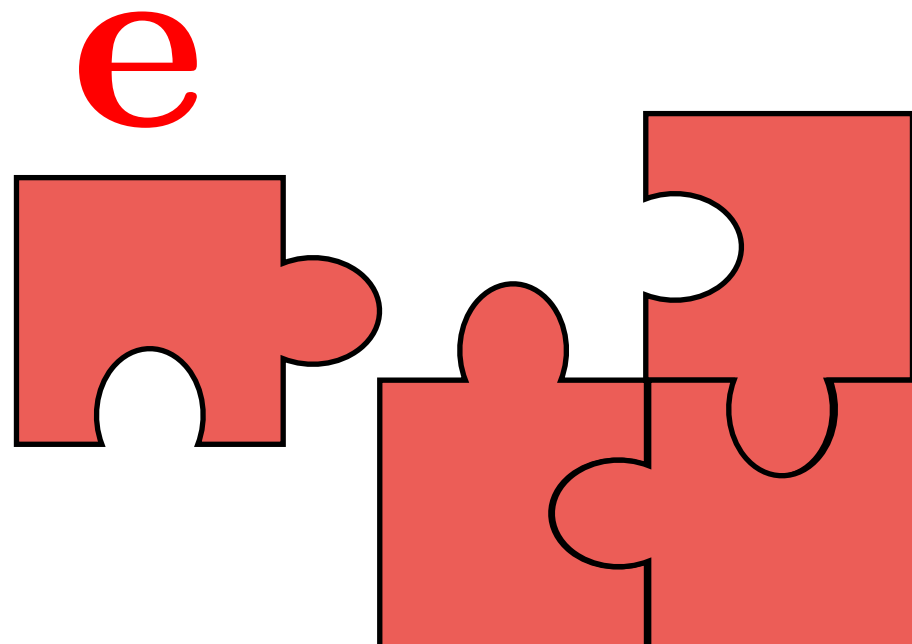
Next...

- Fully Abstract Compilation / Secure Compilation
 - compiler is equivalence-preserving
 - ensures a compiled component does not interact with any target behavior that is inexpressible in S
 - Recent results on fully abstract compilation
- Do we want to link with behavior inexpressible in S ? Or do we want fully abstract compilers?
 - Answer: we want both!
 - How to get there? Languages should let programmers specify what behavior they want to link with

Verified Compilers for Multi-lang. World



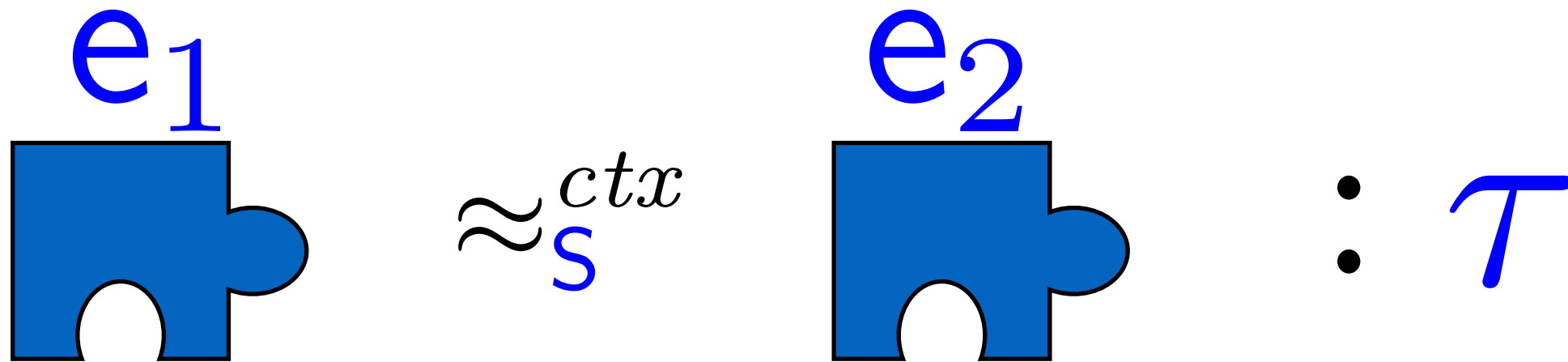
- Runtime
- Drivers



Source programmers should be able to
reason in the **source** language!

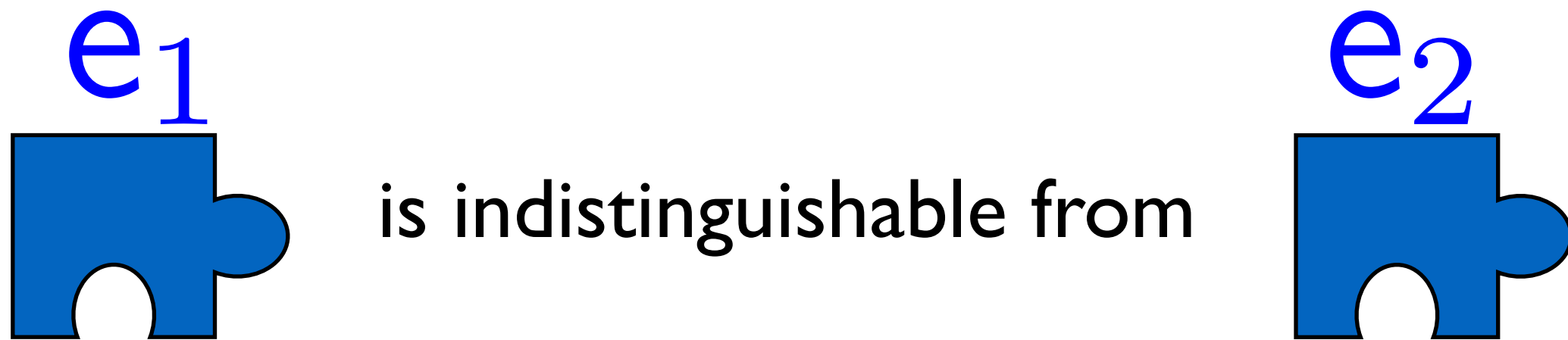
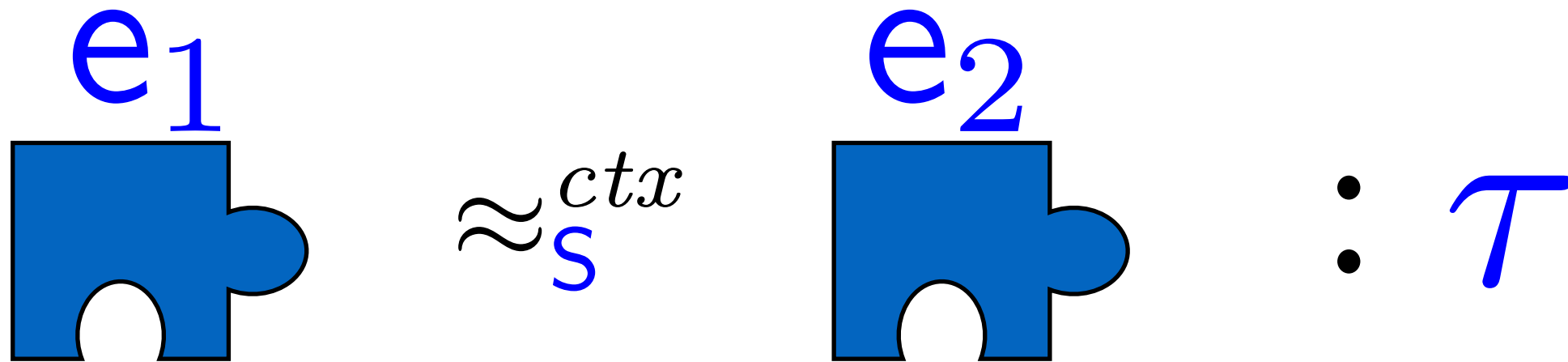
Source Language Reasoning

Contextual Equivalence



Source Language Reasoning

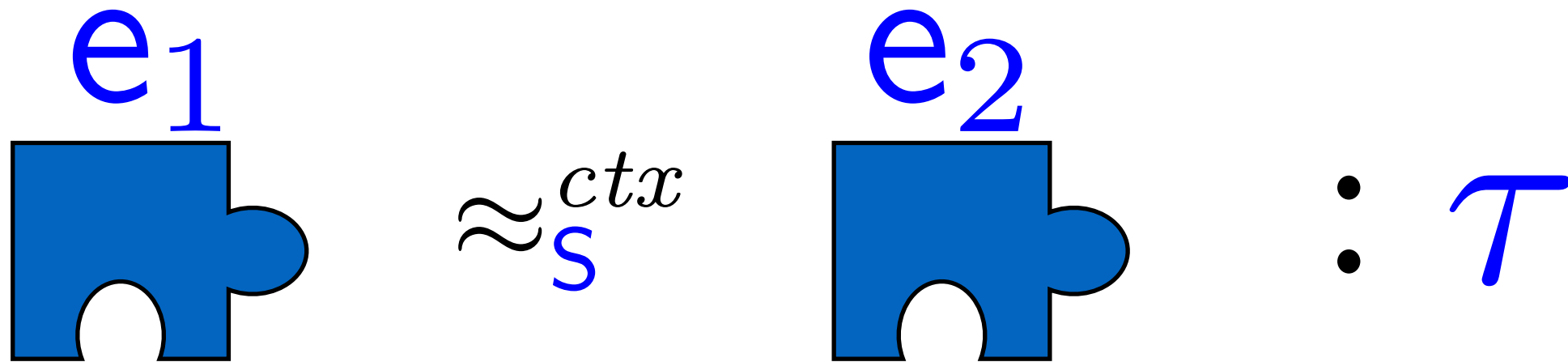
Contextual Equivalence



by contexts at type \mathcal{T}

Source Language Reasoning

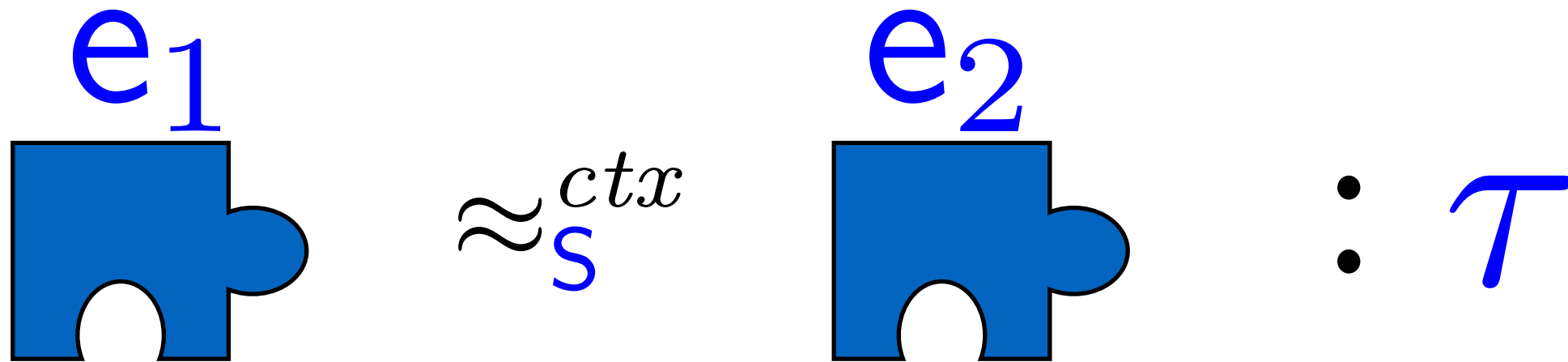
Contextual Equivalence



- Formal basis for
- Refactoring

Source Language Reasoning

Contextual Equivalence

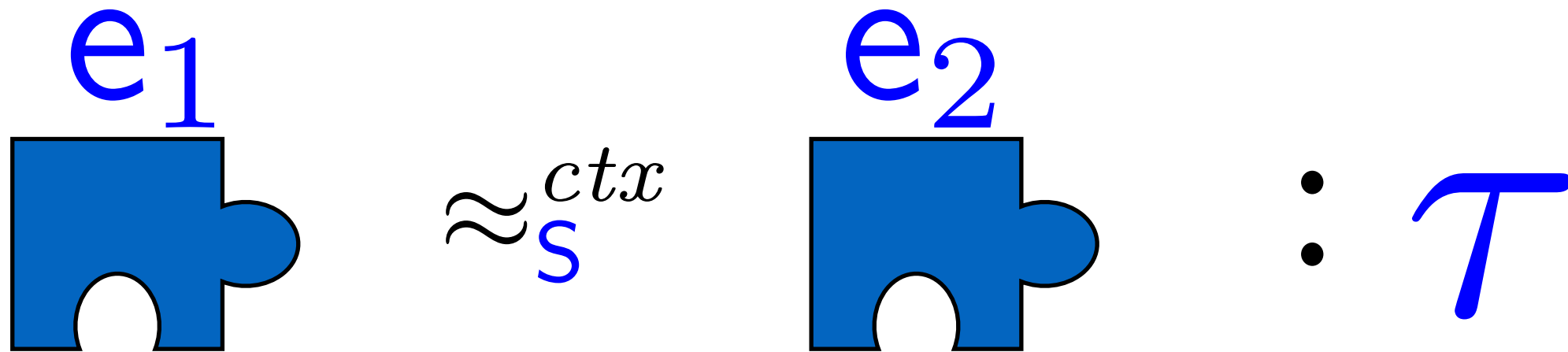


Formal basis for

- Refactoring
- Data abstraction

Source Language Reasoning

Contextual Equivalence

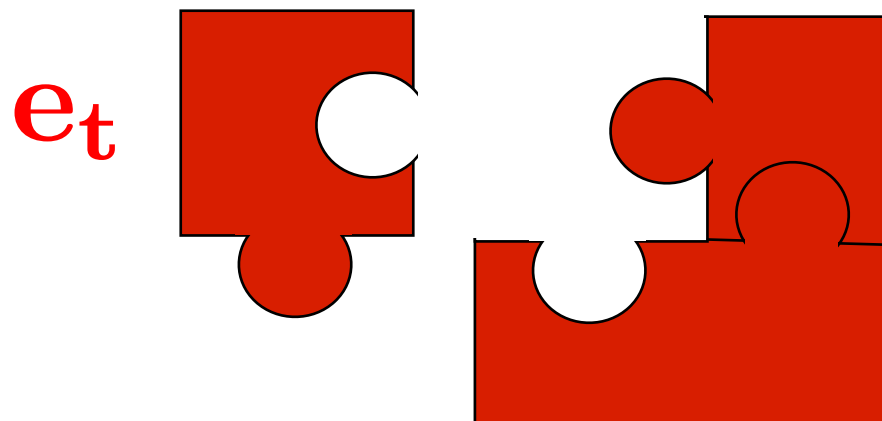
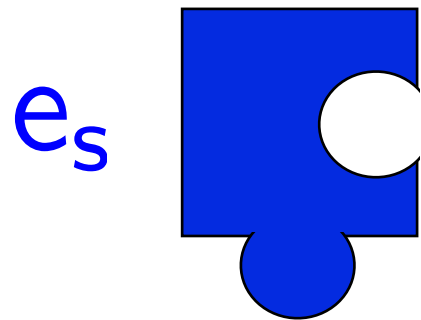


Formal basis for

- Refactoring
- Data abstraction
- Security

Fully Abstract / Secure Compilation

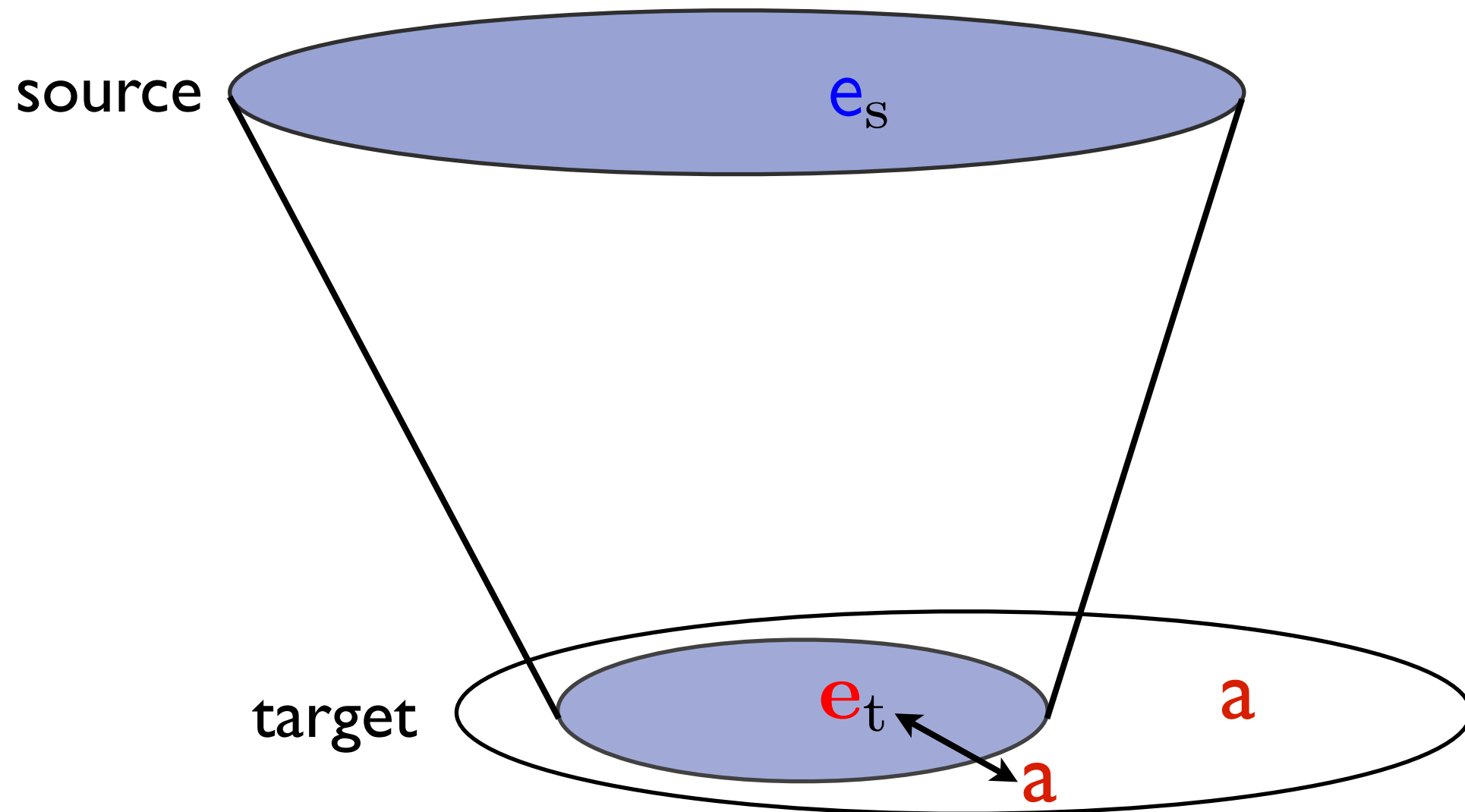
Secure compilation of components:



Want guarantee that e_t will remain as secure as e_s when executed in arbitrary target-level contexts

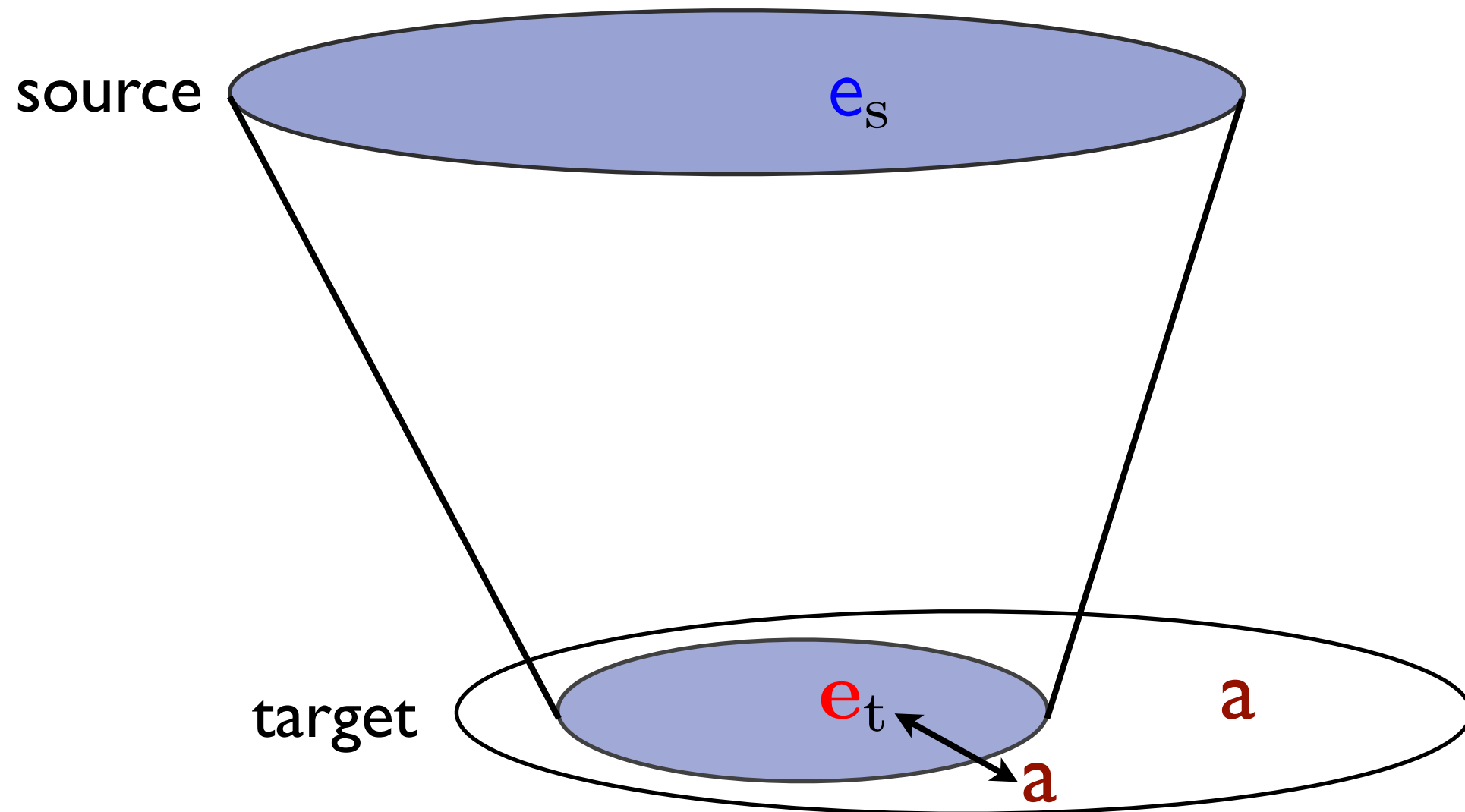
i.e. target contexts (attackers!) can make no more observations about e_t than a source context can make about e_s

Ensuring Secure Compilation



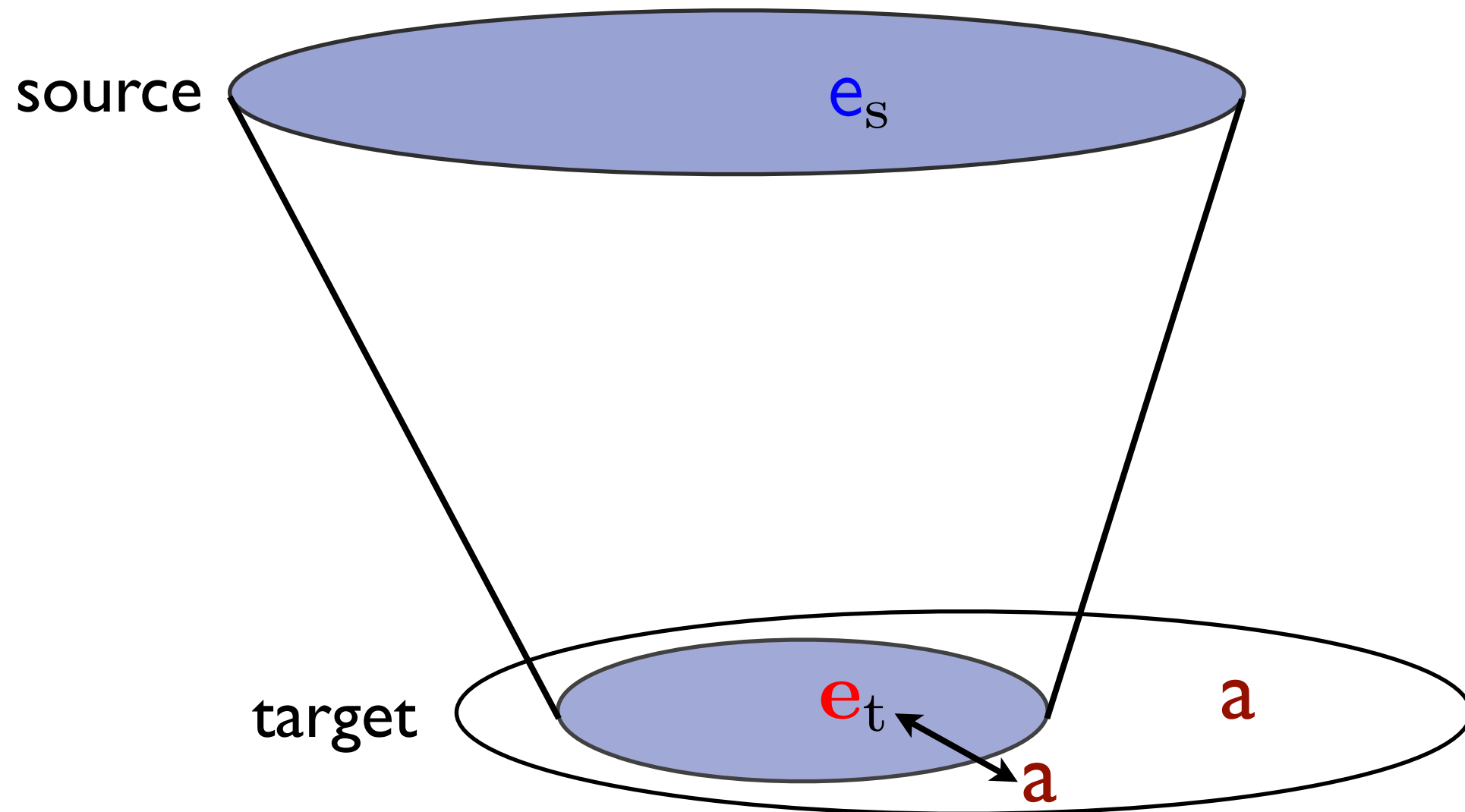
Must ensure that any a we link with behaves like some source context

Ensuring Secure Compilation



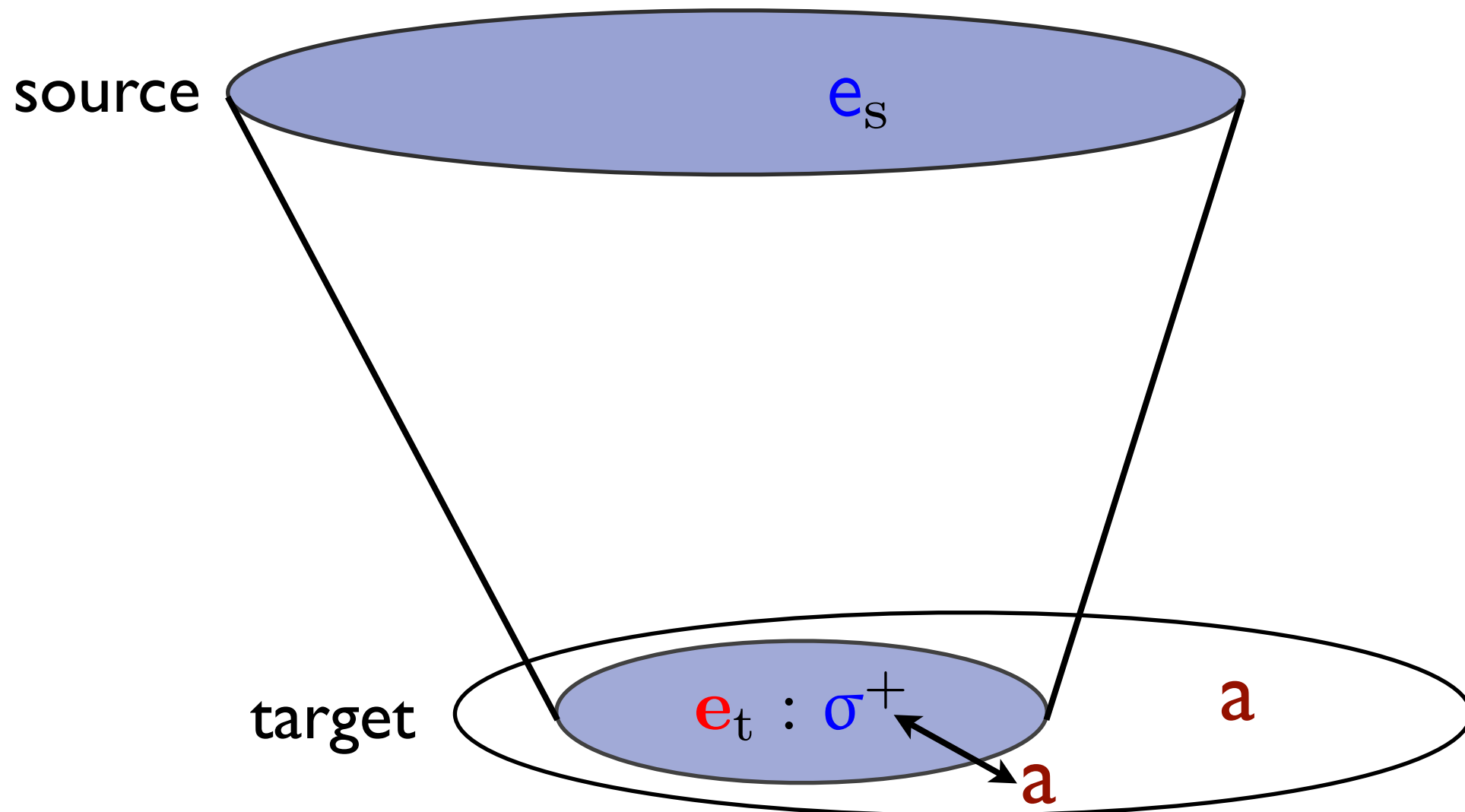
I. Add target features to the source language. **Bad!**

Ensuring Secure Compilation



1. Add target features to the source language. **Bad!**
2. Dynamics checks: catch badly behaved code in the act.
Performance cost

Ensuring Secure Compilation



1. Add target features to the source language. **Bad!**
2. Dynamics checks: catch badly behaved code in the act.
Performance cost
3. Static checks: rule out badly behaved code in the first place
Verification

Type-preserving compilation

$$e : \tau \rightsquigarrow e : \tau^+$$

Equivalence-preserving compilation

If $e_1 : \tau \rightsquigarrow e_1 : \tau^+$ and $e_2 : \tau \rightsquigarrow e_2 : \tau^+$ then:

$$e_1 \approx_S^{ctx} e_2 : \tau \implies e_1 \approx_T^{ctx} e_2 : \tau^+$$

Fully abstract compilation

If $e_1 : \tau \rightsquigarrow e_1 : \tau^+$ and $e_2 : \tau \rightsquigarrow e_2 : \tau^+$ then:

$$e_1 \approx_S^{ctx} e_2 : \tau \iff e_1 \approx_T^{ctx} e_2 : \tau^+$$

preserves & reflects equivalence



Challenge of **proving** full abstraction

Suppose $\Gamma \vdash e_1 : \tau \rightsquigarrow e_1$ and $\Gamma \vdash e_2 : \tau \rightsquigarrow e_2$.

$$\Gamma \vdash e_1 \approx_S^{ctx} e_2 : \tau$$



$$\Gamma^+ \vdash e_1 \approx_T^{ctx} e_2 : \tau^+$$

Challenge of **proving** full abstraction

Suppose $\Gamma \vdash e_1 : \tau \rightsquigarrow e_1$ and $\Gamma \vdash e_2 : \tau \rightsquigarrow e_2$.

$$\Gamma \vdash e_1 \approx_S^{ctx} e_2 : \tau$$



$$\Gamma^+ \vdash e_1 \approx_T^{ctx} e_2 : \tau^+$$

Given:

No C_S can
distinguish e_1, e_2

Challenge: Back-translation

- I. **If target is not more expressive than source**, use the same language: back-translation can be avoided in lieu of *wrappers* between τ and τ^+
 - Closure conversion: System F with recursive types
[Ahmed-Blume ICFP'08]
 - f^* (STLC with refs) to js^* (encoding of JavaScript in f^*)
[Fournet et al. POPL'13]

Challenge: Back-translation

2. If target is more expressive than source

(a) Both **terminating**: use back-translation by partial evaluation

- Equivalence-preserving CPS from STLC to System F
[Ahmed-Blume ICFP'11]
- Noninterference for Free (DCC to F_ω)
[Bowman-Ahmed ICFP'15]

(b) Both **nonterminating**: use ??

back-trans by partial evaluation is not well-founded!

Observation: our source lang. has recursive types,
can write interpreter for target lang. in source lang.

Fully Abstract Closure Conversion

[New et al.'16]

Source: STLC + μ types

Target: System F + \exists types + μ types + **exceptions**

First full abstraction result where target has exceptions but source does not.

Earlier work, due to lack of sufficiently powerful back-translation techniques, adds features from target to source.

Novel proof technique — **Universal Embedding**

- Untyped embedding of target in source
- Mediate between strongly typed source and untyped back-translation

Fully Abstract Closure Conversion

Source: STLC + μ types

Target: System F + \exists types + μ types + **exceptions**

Equivalent source terms, inequivalent in lang. with exceptions:

$e_1 = \lambda f. (f \text{ true}; f \text{ false}; \langle \rangle)$ $e_2 = \lambda f. (f \text{ false}; f \text{ true}; \langle \rangle)$

$C = \text{catch } y = ([\cdot] (\lambda x. \text{raise } x)) \text{ in } y$

$C[e_1] \Downarrow \text{true}$

$C[e_2] \Downarrow \text{false}$

Idea: use modal type system at target to rule out linking with code that throws unhandled exceptions

Ensuring Full Abstraction

$$e_1 \approx_S^{ctx} e_2 : (\text{bool} \rightarrow 1) \rightarrow 1$$

$$(\text{bool} \rightarrow \mathbf{E} 0 1) \rightarrow \mathbf{E} 0 1$$

\neq

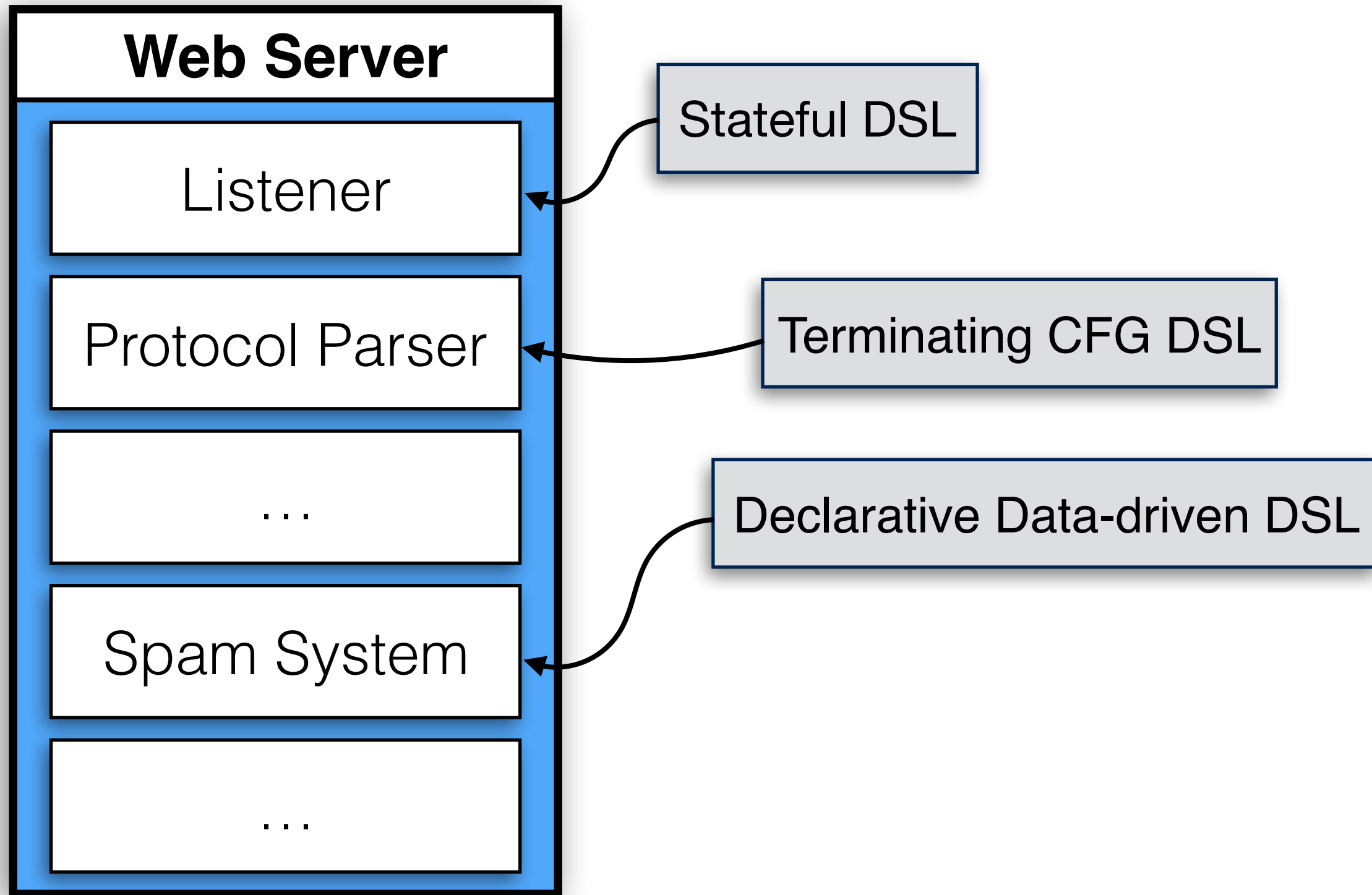
$$\mathbf{C} : (\text{bool} \rightarrow \mathbf{E} \text{ bool } 1) \rightarrow \mathbf{E} \text{ bool } 1$$

$$\mathbf{C} = ([\cdot] \lambda(x : \text{bool}). \text{raise } x)$$

Static Fully Abstract Compilation

- Type checking ensures that we never link with target code whose (extensional) behavior does not match *some* source behavior
- But what if we want to link with behaviors unavailable in the source?
 - Surely, we want that when building multi-language software!

Multi-Language System



Web Server

Protocol Parser

Terminating DSL

Listener

Stateful DSL

???

Does Nontermination Leak?

COMPILE

Can we instead allow reasoning at source level?

COMPILE

Protocol Parser

Common Target

Listener

Common Target

Linking

Linking types are about raising programmer reasoning back to the source level

Linking Types for Multi-Language Software:
Have Your Cake and Eat it Too
[Patterson-Ahmed SNAPL'17]

In a Simpler Setting

λ (simply-typed
lambda calculus)

$\tau ::= \mathbf{unit} \mid \mathbf{int} \mid \tau \rightarrow \tau$
 $e ::= () \mid \mathbf{n} \mid \mathbf{x} \mid \lambda \mathbf{x} : \tau. e$
 $\mid ee \mid e + e \mid e * e$

λ^{ref} (extended with
ML references)

$\tau ::= \dots \mid \text{ref } \tau$
 $e ::= \dots \mid \text{ref } e \mid e := e \mid !e$

How to reason in λ while linking with λ^{ref} ?

Refactoring is reasoning about equivalence

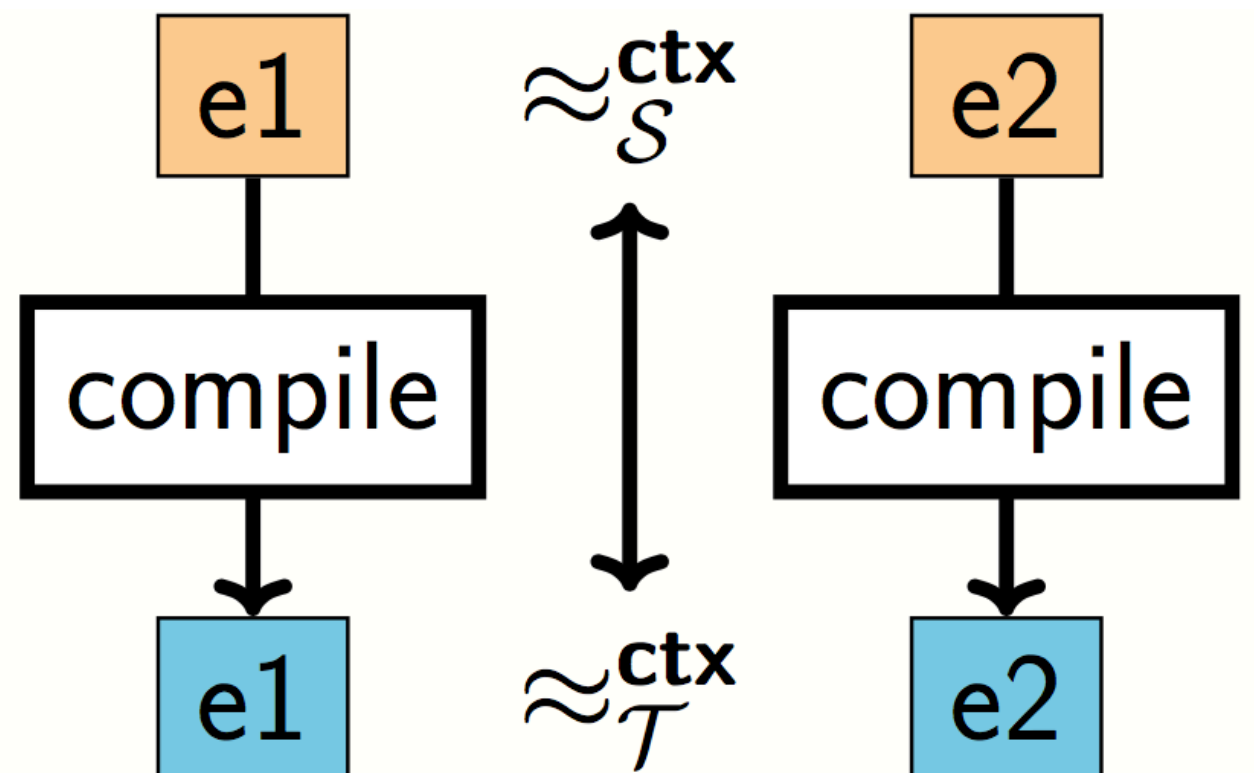
Reasoning About Refactoring

$\lambda c. c(); c() \Rightarrow \lambda c. c() : (\text{unit} \rightarrow \text{int}) \rightarrow \text{int}$

Should be okay because

$\lambda c. c(); c() \approx_{\lambda}^{ctx} \lambda c. c()$

Fully abstract
compilers preserve
equivalences



What about linking with λ^{ref} ?

```
let counter f' = let v = ref 0 in  
    let c' () = v := !v + 1; !v in f' c'  
let f  
in counter f
```

⇓ 2

but

```
let counter f' = let v = ref 0 in  
    let c' () = v := !v + 1; !v in f' c'  
let f  
in counter f
```

⇓ 1

When linked with λ^{ref} , no longer equivalent!

Is this refactoring correct?

$\lambda c. c(); c()$ \Rightarrow $\lambda c. c() : (\text{unit} \rightarrow \text{int}) \rightarrow \text{int}$

It depends on what it is linked with!

$\text{unit} \rightarrow \text{int}$



$\text{unit} \rightarrow \text{int}$



Programmer should be able to specify which they want, so that the compiler can be fully abstract!

λ with linking types extension

$$\tau ::= \text{unit} \mid \text{int} \mid \tau \rightarrow \tau$$

$$\lambda^{\kappa} \left| \begin{array}{l} \tau ::= \text{unit} \mid \text{int} \mid \tau \rightarrow \mathbf{R}^{\emptyset} \tau \\ \mid \text{ref } \tau \mid \tau \rightarrow \mathbf{R}^{\downarrow} \tau \end{array} \right.$$

Type and effect systems, e.g., F^* , Koka

λ^{κ} Allows Programmers To Write Both

unit \rightarrow **int**

unit \rightarrow R^{\emptyset} **int**

unit \rightarrow **int**

unit \rightarrow R^{\downarrow} **int**

Refactoring: Pure Inputs

$\lambda c: \text{unit} \rightarrow R^\emptyset \text{int}. c(); c() \approx_{\lambda \kappa}^{\text{ctx}} \lambda c: \text{unit} \rightarrow R^\emptyset \text{int}. c()$

```
let counter f' = let v = ref 0 in  
  let c' () = v := !v + 1; !v in f' c'  
let f = ... =  $\lambda c: \text{unit} \rightarrow R^\emptyset \text{int}. c()$   
in counter f
```

Ill-typed, since `f` requires pure code

Refactoring: Impure Inputs

$\lambda c: \text{unit} \rightarrow \mathbf{R}^{\downarrow} \text{int}. c(); c() \not\approx_{\lambda \kappa}^{ctx} \lambda c: \text{unit} \rightarrow \mathbf{R}^{\downarrow} \text{int}. c()$

let counter $f' = \text{let } v = \text{ref } 0 \text{ in}$
 let $c' () = v := !v + 1; !v \text{ in } f' c'$
let $f = \lambda c: \text{unit} \rightarrow \mathbf{R}^{\downarrow} \text{int}. c()$
in counter f

Well-typed, since f accepts impure code

Minimal Annotation Burden

$$\lambda c: \text{unit} \rightarrow \mathbf{R}^\emptyset \text{int}. c(); c()$$

$$\lambda c: \text{unit} \rightarrow \text{int}. c(); c()$$

λ^κ must provide default translation

$$\kappa^+(\text{unit}) = \text{unit}$$

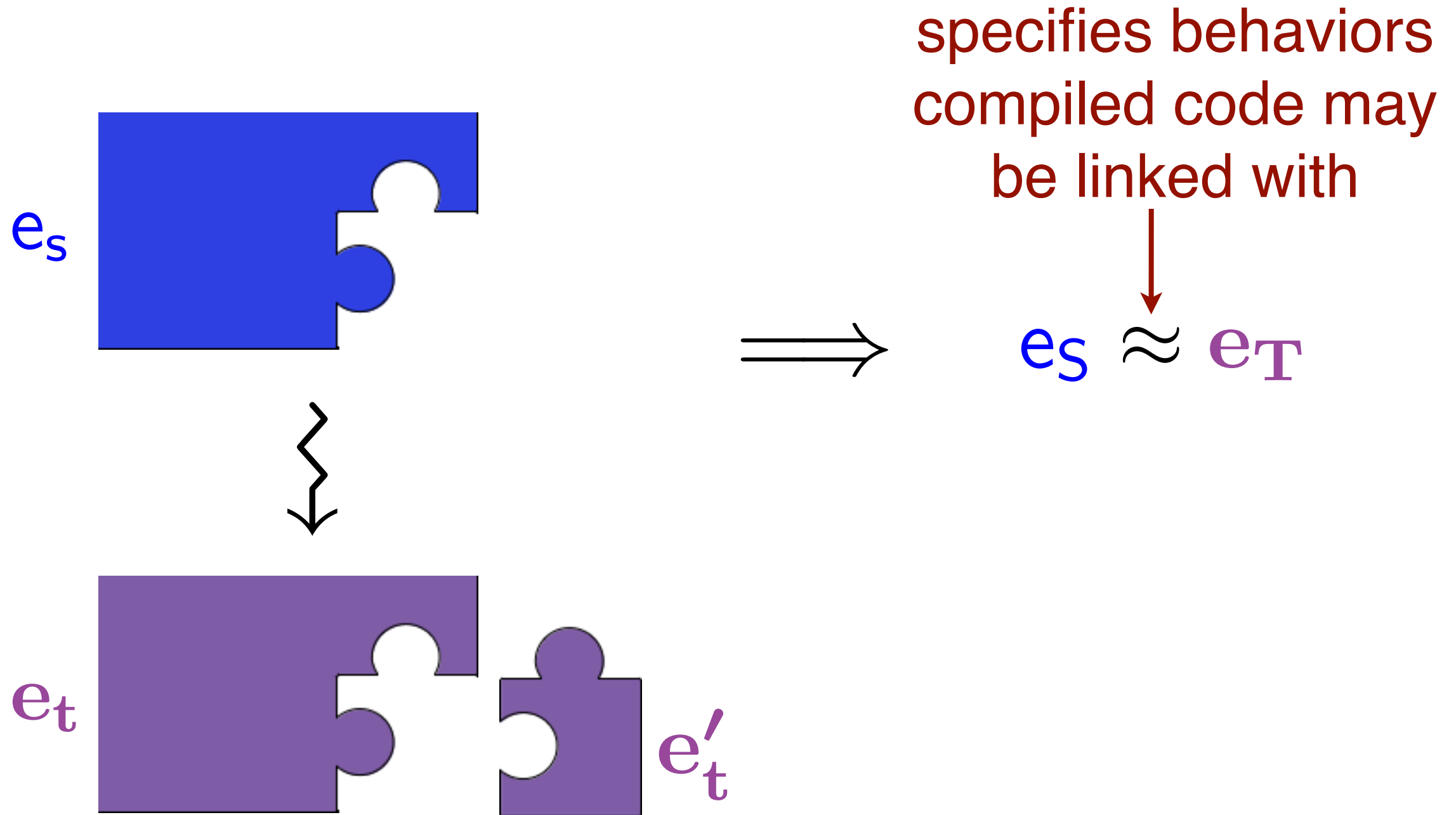
$$\kappa^+(\text{int}) = \text{int}$$

$$\kappa^+(\tau_1 \rightarrow \tau_2) = \kappa^+(\tau_1) \rightarrow \mathbf{R}^\emptyset \kappa^+(\tau_2)$$

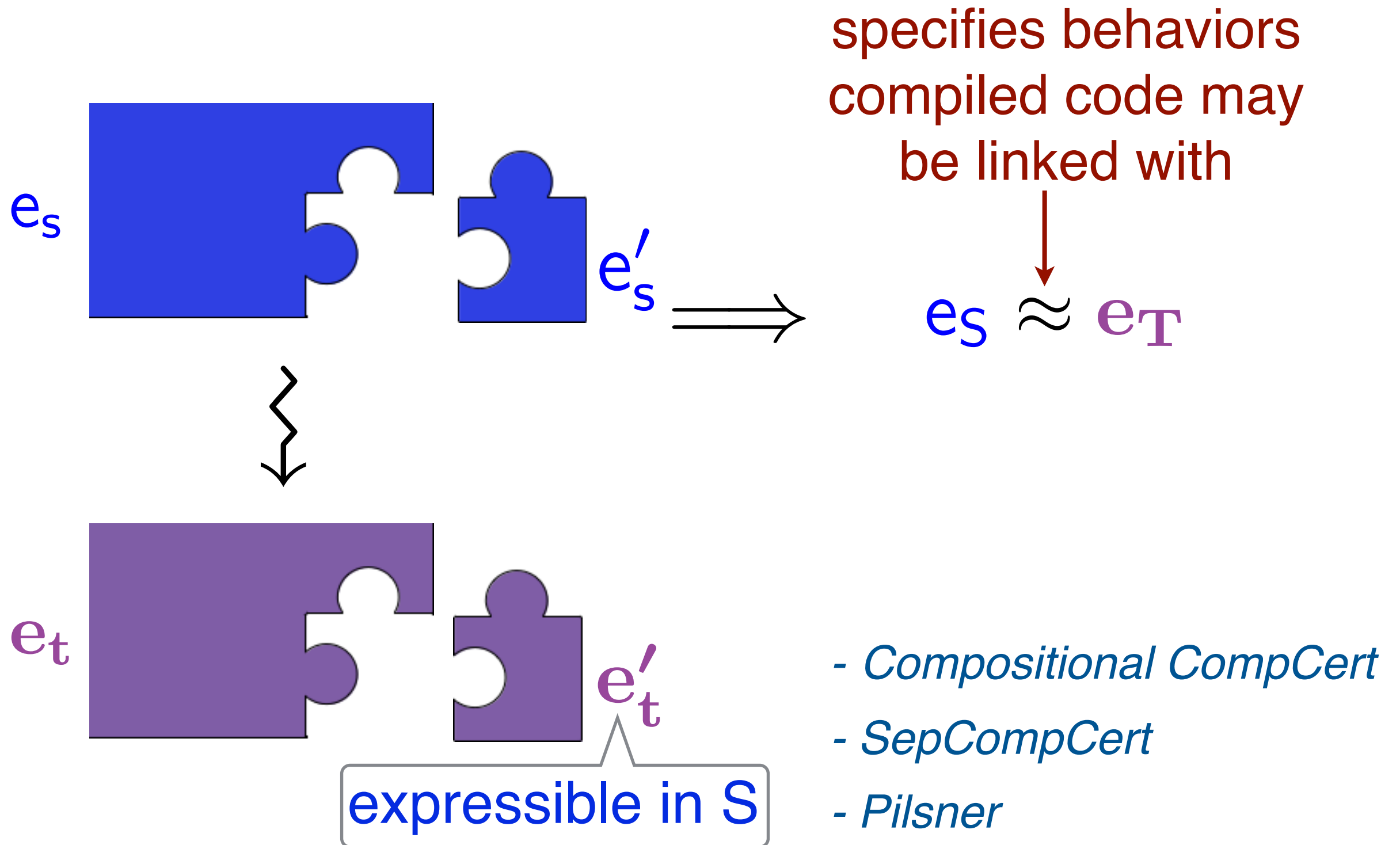
$$\forall e_1, e_2. e_1 \approx_{\lambda}^{\text{ctx}} e_2 : \tau \implies e_1 \approx_{\lambda^\kappa}^{\text{ctx}} e_2 : \kappa^+(\tau)$$

Stepping Back...

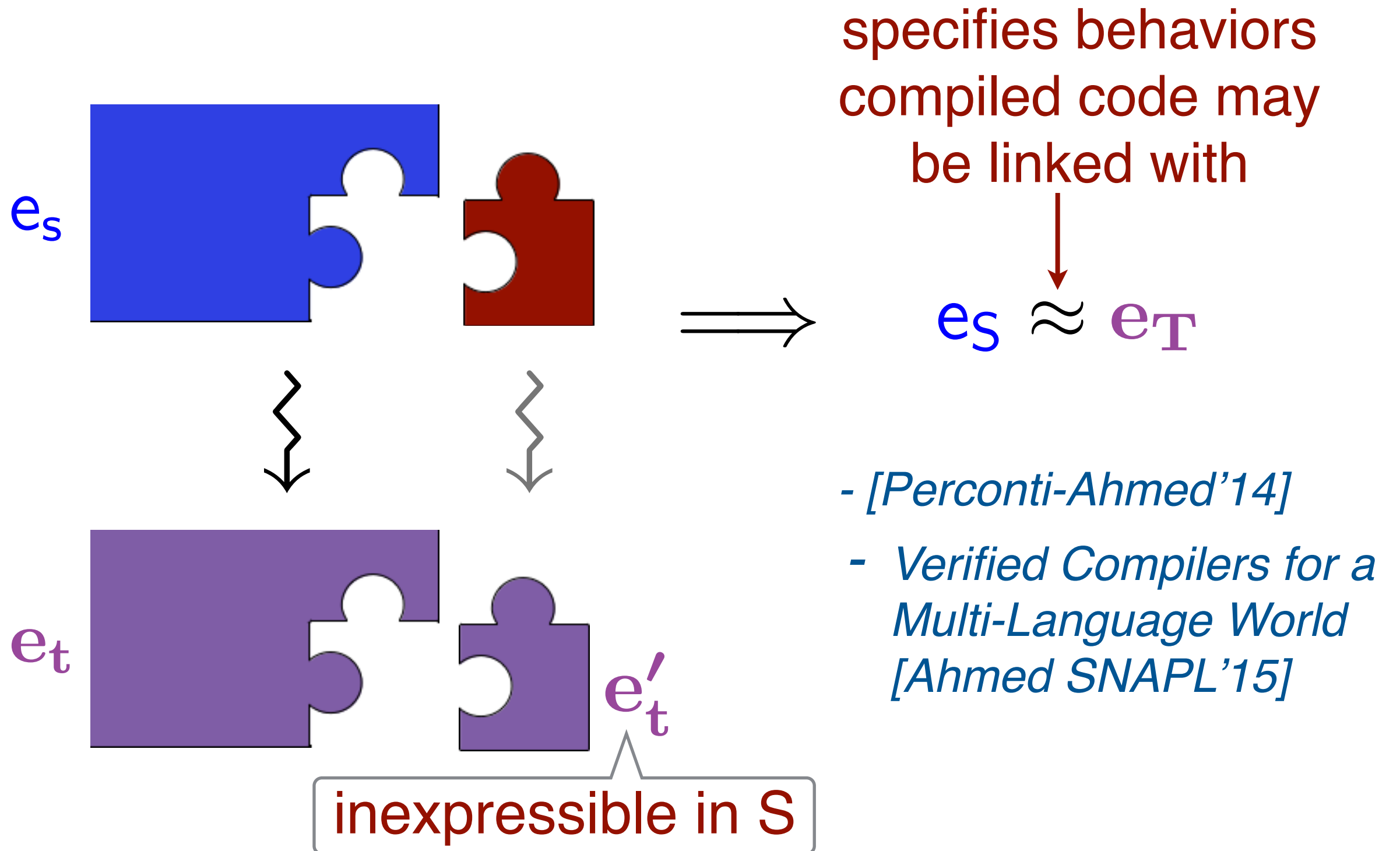
Correct Compilation of Components



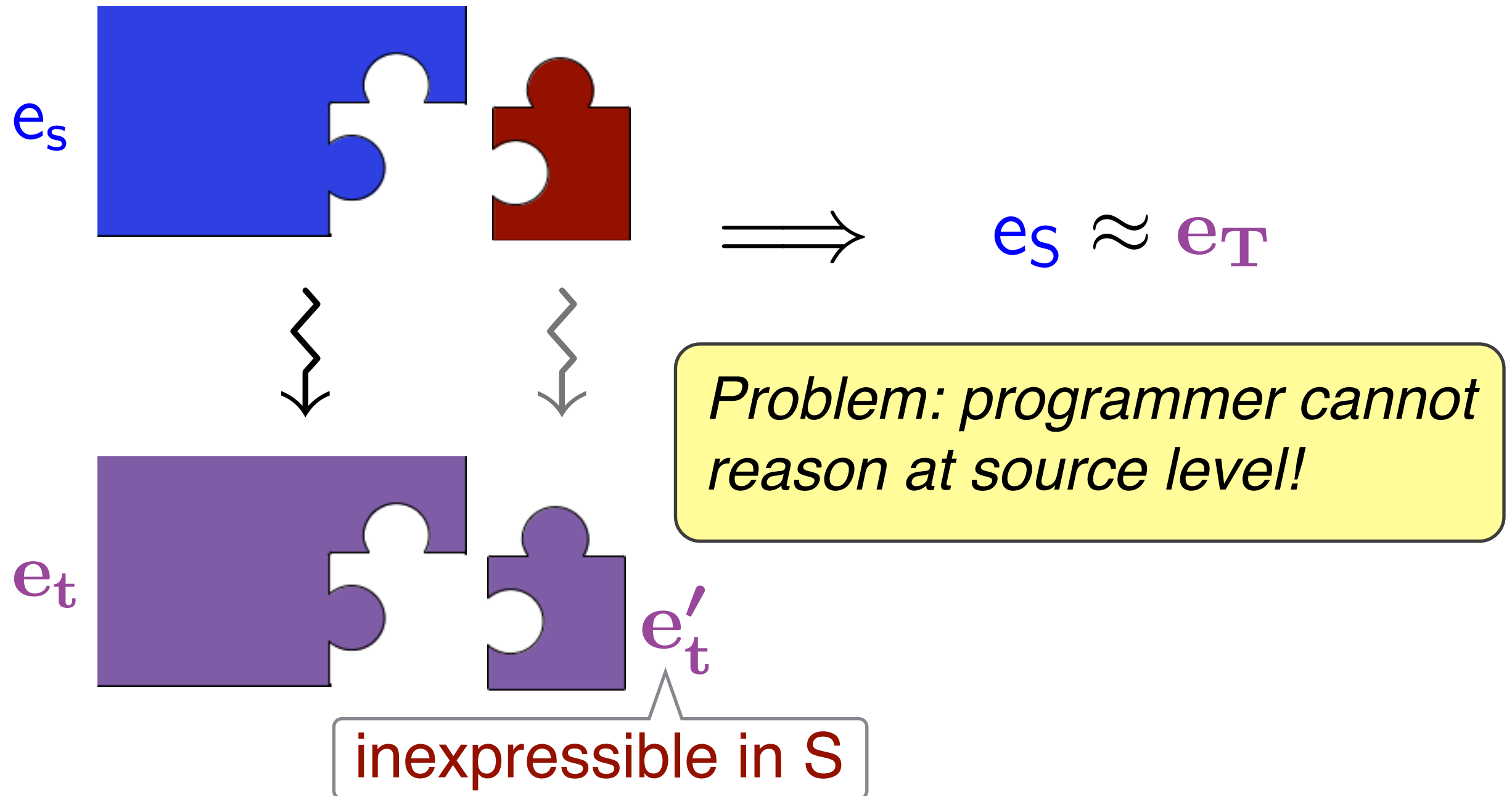
Correct Compilation of Components



Correct Compilation: Multi-Language



Correct Compilation: Multi-Language



Fully Abstract Compilation?

*escape
hatches*

ML
C FFI

Rust
unsafe

Java
JNI

*Language specifications are incomplete!
Don't account for linking*

Target

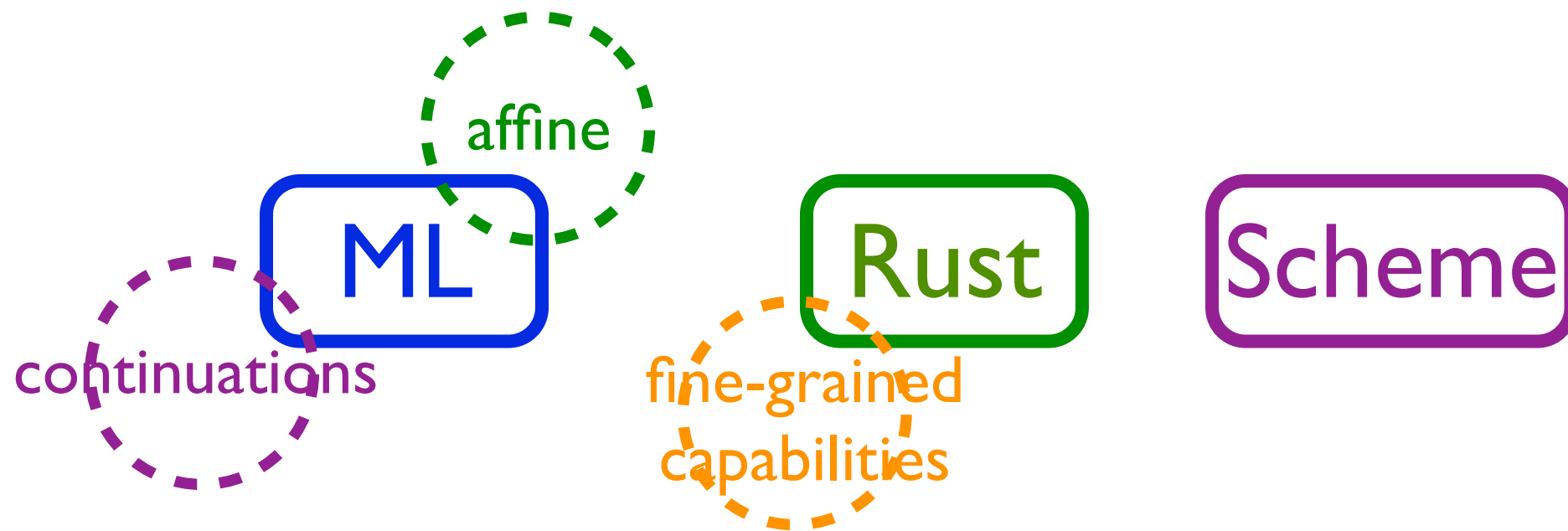
Rethink PL Design with Linking Types

*escape
hatches*



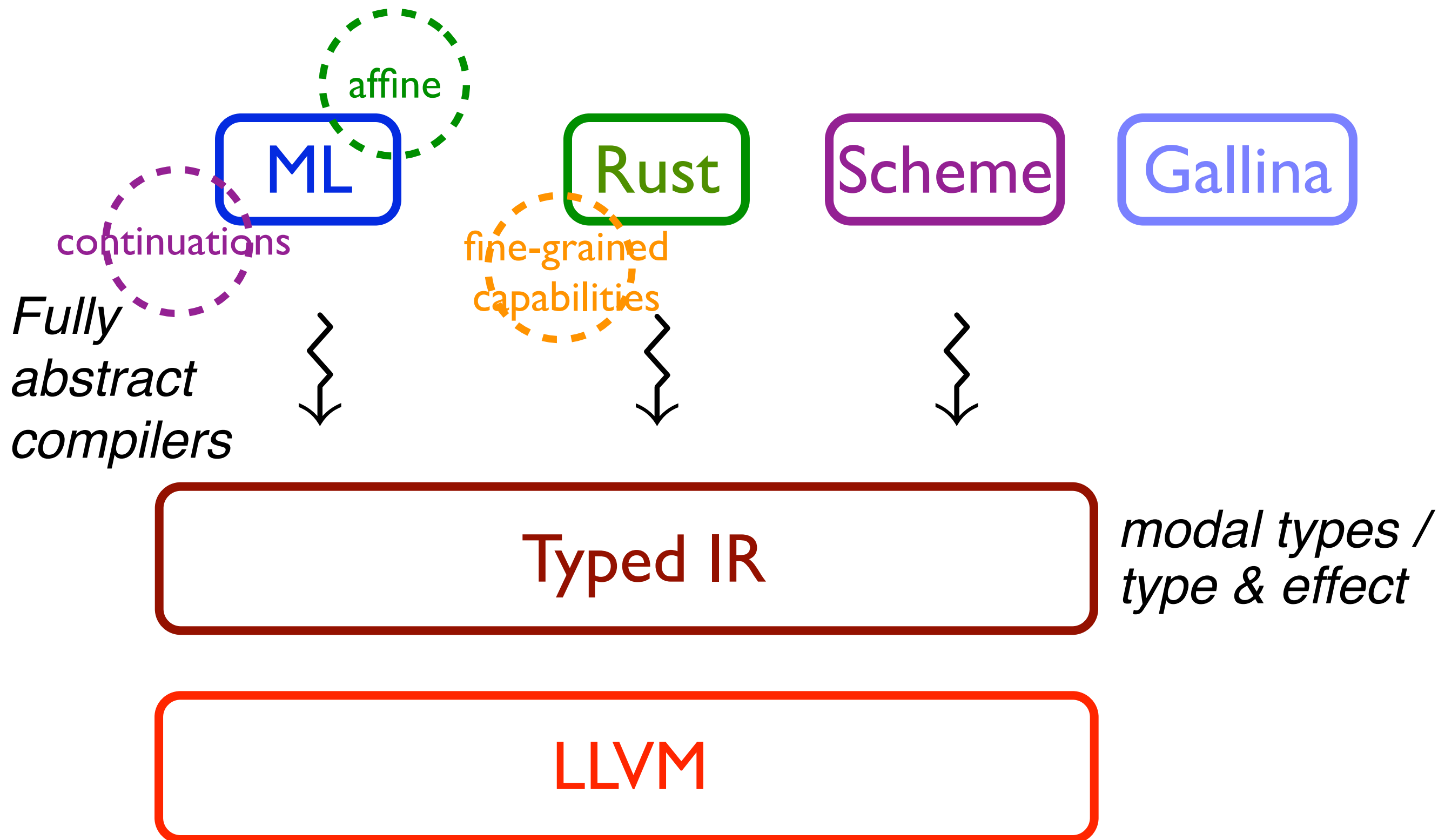
Design linking types extensions that support safe interoperability with other languages

PL Design, Linking Types

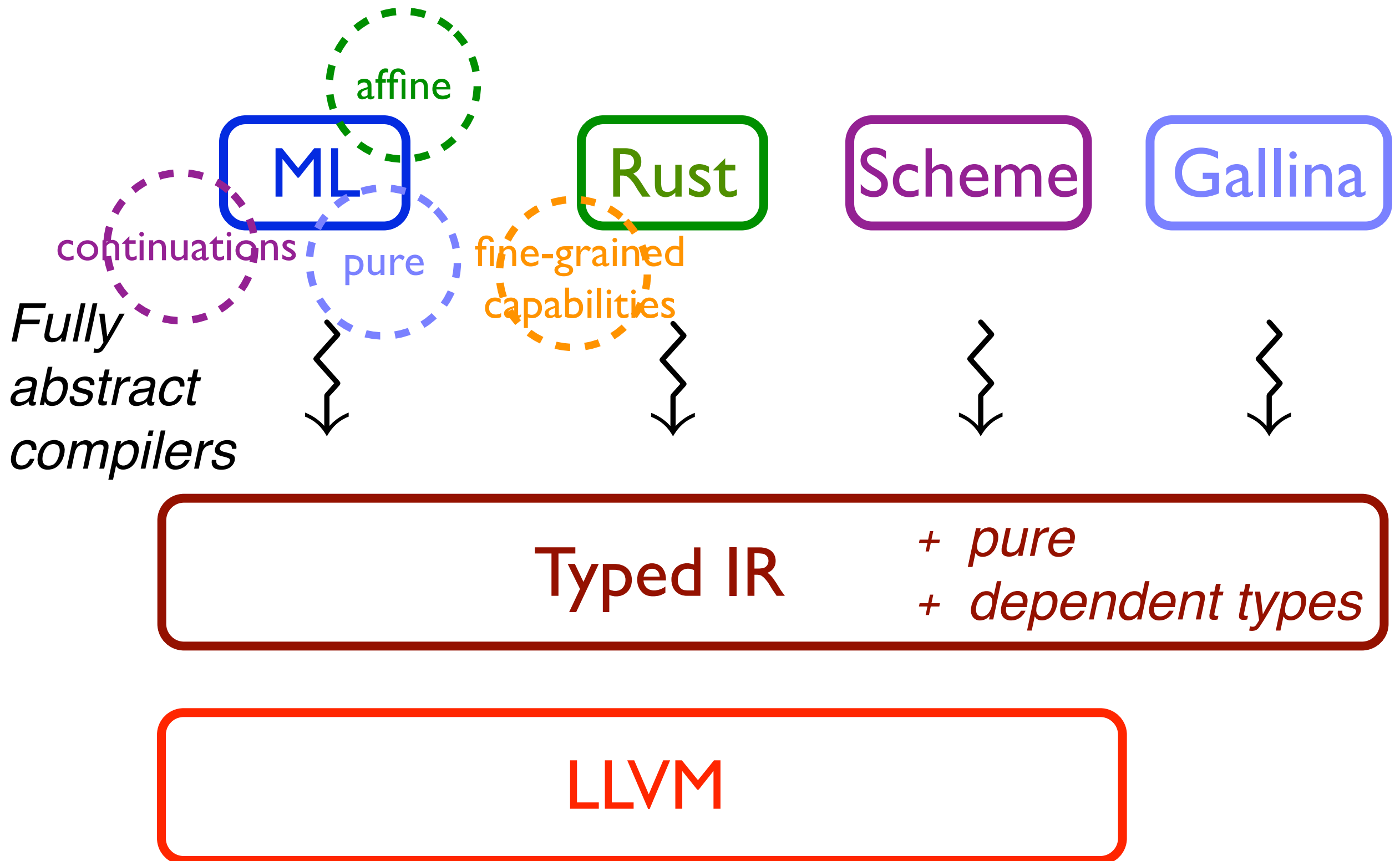


Only need linking types extensions to interact with behavior inexpressible in your language.

PL Design, Linking Types, Compilers



PL Design, Linking Types, Compilers



Linking Types

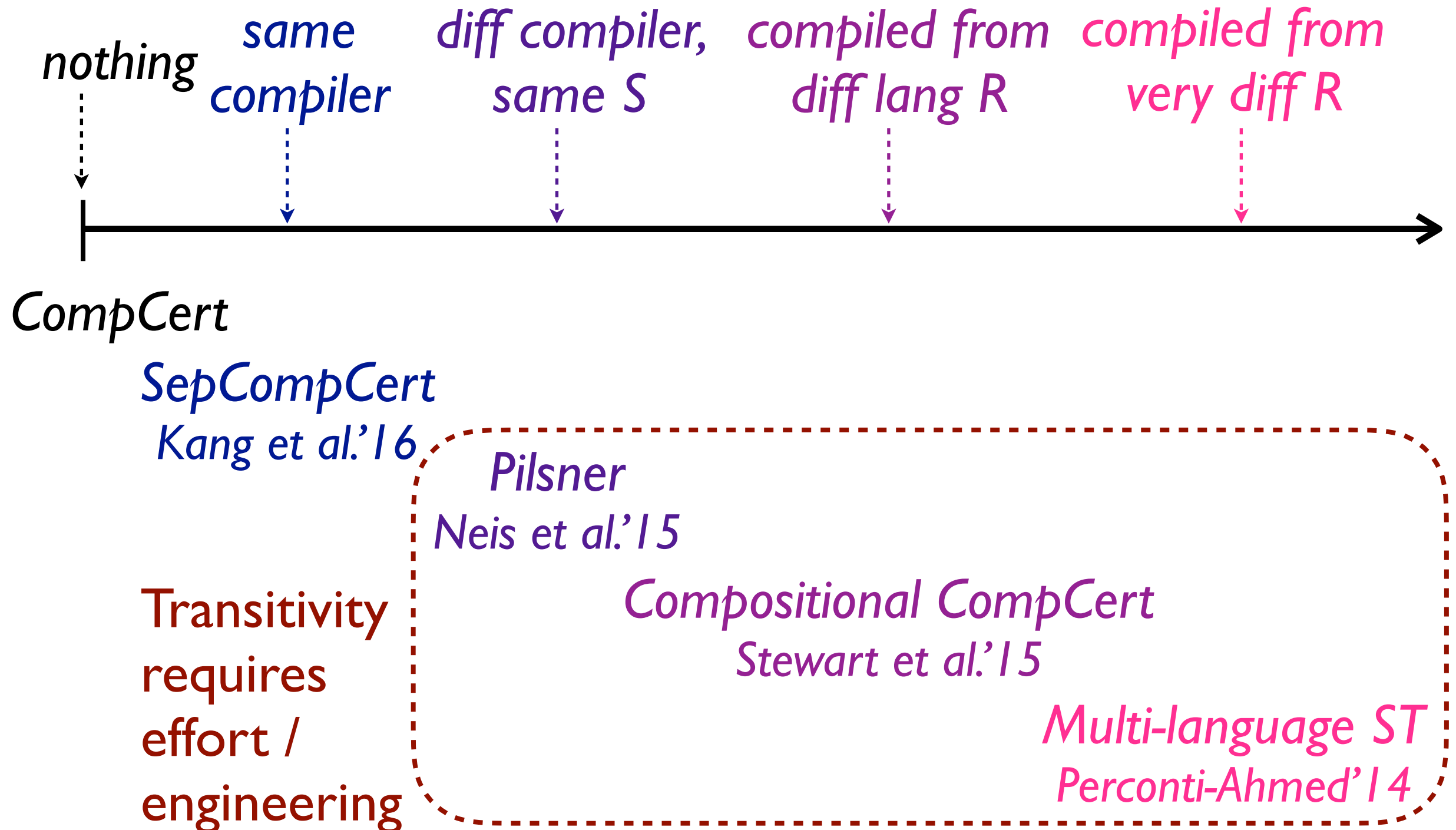
- Allow programmers to reason in *almost* their own source languages, even when building multi-language software
- Allow compilers to be fully abstract, yet support multi-language linking

Conclusion

Compiler Verif. for Multi-Lang. World

- Compositional Compiler Correctness
 - horizontal and vertical compositionality

Horizontal / **Vertical** Compositionality



Compiler Verif. for Multi-Lang. World

- CompCert started a renaissance in compiler verification
 - major advances in mechanized proof
- Now we need: Compositional Compiler Correctness
 - but horizontal and vertical compositionality at odds
- Need to rethink proof architectures for compiler verification to support linking with code of arbitrary provenance. But want transitivity to be easier!

Verification of *realistic* compilers for a multi-language world demands formal techniques and language design

- **compositional equational reasoning**
- formal semantics of **language interoperability**
- types and logics to enforce **sensible (safe, secure) linking**
- extending our language designs with *principled extensions* to replace unprincipled escape hatches