



Manuel Di Agostino

`manuel.diagostino@studenti.unipr.it`

Università di Parma

December 10, 2025

A generic framework for heap and value analyses of object-oriented programming languages

Based on: Ferrara, P. (2016). Theoretical Computer Science, 631, 43-72.

1 Introduction

- The problem

2 The framework

- Concrete domain and semantics
- Abstract domain and semantics
- Instantiation

3 Conclusion and Contributions

1 Introduction

- The problem

2 The framework

- Concrete domain and semantics
- Abstract domain and semantics
- Instantiation

3 Conclusion and Contributions

We would like to say something about the following:

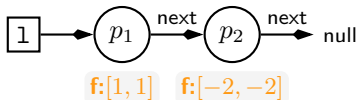
```
1 int absSum(ListInt l) {  
2   int sum = 0;  
3   ListInt it = l;  
4  
5   while(it != null) {  
6     if (it.f < 0)  
7       it.f = -it.f;  
8  
9     sum += it.f  
10    it = it.next;  
11  }  
12  return sum;  
13 }
```

***Listing:** Sum of absolute values of a list of integers.*

Two clients, different needs:

➤ **Client 1:** `absSum([1, -2])`

➤➤ Fixed list, 2 nodes at labels p_1 , p_2



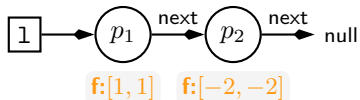
```
1 int absSum(ListInt l) {
2     int sum = 0;
3     ListInt it = l;
4
5     while(it != null) {
6         if (it.f < 0)
7             it.f = -it.f;
8
9         sum += it.f
10        it = it.next;
11    }
12    return sum;
13 }
```

The problem

Two clients, different needs:

➤ **Client 1:** `absSum([1, -2])`

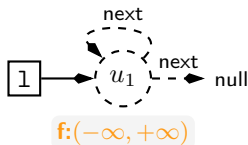
➤➤ Fixed list, 2 nodes at labels p_1, p_2



➤ **Client 2:** `absSum($[l_i \in \mathbb{Z}]$)`,
 $1 \leq i \leq n$

➤➤ Dynamic list, n unknown

➤➤ Summary node u_1 abstracts all nodes

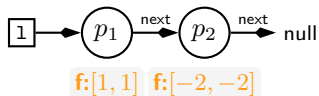


```

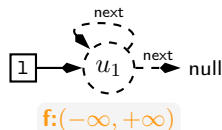
1 int absSum(ListInt l) {
2   int sum = 0;
3   ListInt it = l;
4
5   while(it != null) {
6     if (it.f < 0)
7       it.f = -it.f;
8
9     sum += it.f
10    it = it.next;
11  }
12  return sum;
13 }
```

The Allocation-Site method

Client 1 (fixed list)



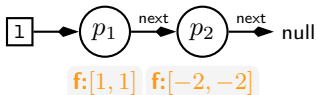
Client 2 (dynamic list)



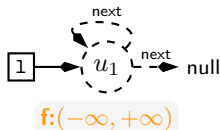
Property	Client 1 Result	Client 2 Result
P1: No NullPointer (Heap Structure)	✓ Safe Heap structure is precise	✓ Safe Summary handle nulls

The Allocation-Site method

Client 1 (fixed list)



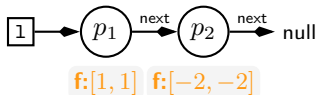
Client 2 (dynamic list)



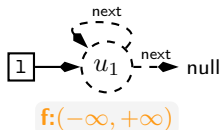
Property	Client 1 Result	Client 2 Result
P1: No NullPointer (Heap Structure)	✓ Safe Heap structure is precise	✓ Safe Summary handle nulls
P2: Return ≥ 0 (Sign Analysis)	✓ Verified Strong updates on p_1, p_2	✗ False Alarm Weak update on u_1

The Allocation-Site method

Client 1 (fixed list)



Client 2 (dynamic list)



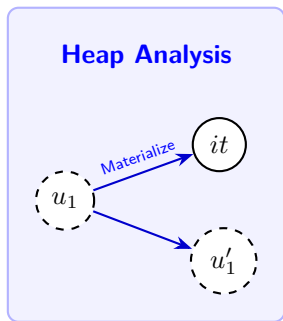
Property	Client 1 Result	Client 2 Result
P1: No NullPointer (Heap Structure)	✓ Safe Heap structure is precise	✓ Safe Summary handle nulls
P2: Return ≥ 0 (Sign Analysis)	✓ Verified Strong updates on p_1, p_2	✗ False Alarm Weak update on u_1
P3: Elem ≥ 0 (Relational)	✓ Verified Relations preserved	✗ False Alarm Relations lost on u_1

Materialization: The First Step

The Shape Analysis

When iterating the loop, the Heap Analysis **must materialize** the summary node u_1 to distinguish the current element from the rest:

- A concrete node it (the current element).
- A remaining summary node u'_1 (the rest of the list).



NO AUTOMATIC
COMMUNICATION

Value Analysis

Pre-state:

$$u_1.f \in (-\infty, +\infty)$$

Post-state:

$$it.f \in ?$$

$$u'_1.f \in ?$$

The Consequence: Information Loss

What we miss: When u_1 splits into it and u'_1 , where should the information $u_1.f \in (-\infty, +\infty)$ go?

Before		After Materialization
$u_1.f \in (-\infty, +\infty)$	$\xrightarrow{\text{materialize}}$	$\begin{cases} it.f \in (-\infty, +\infty) \\ u'_1.f \in (-\infty, +\infty) \end{cases}$

Remark (Consequence: False Alarms)

Without guidance from the Heap Analysis, the Value Domain conservatively assumes **top** for both $it.f$ and $u'_1.f$. This forces analysis of the branch 'if ($it.f < 0$)', causing false alarms on properties P2 and P3.

Key Idea: The Heap Analysis communicates how identifiers are transformed through a **substitution** message:

$$\{it.f, u'_1.f\} \mapsto \{u_1.f\}$$

What the framework does:

- 1 Heap Analysis materializes $\{it, u'_1\} \mapsto \{u_1\}$
- 2 Produces substitution: $\{it.f, u'_1.f\} \mapsto \{u_1.f\}$
- 3 Value Domain receives substitution
- 4 Value Domain **automatically updates** its state:
 - » $it.f \in (-\infty, +\infty)$
 - » $u'_1.f \in (-\infty, +\infty)$

Key Idea: The Heap Analysis communicates how identifiers are transformed through a **substitution** message:

$$\{it.f, u'_1.f\} \mapsto \{u_1.f\}$$

Remark

- Inside the branch `if (it.f < 0)`, the analysis knows *it.f* is negative.
- Since *it* is materialized, `it.f = -it.f` performs a **strong update**.
- The value is *overwritten* to $(0, +\infty)$, so we sum only **non-negative values**!

1 Introduction

- The problem

2 The framework

- Concrete domain and semantics
- Abstract domain and semantics
- Instantiation

3 Conclusion and Contributions

1 Introduction

- The problem

2 The framework

- Concrete domain and semantics
- Abstract domain and semantics
- Instantiation

3 Conclusion and Contributions



In a concrete execution, a standard object-oriented state holds two distinct types of data:

- **Values** (Val): Primitive data (integers, booleans, etc.)
- **References** (Ref): Pointers to objects or `null`



Information

Common in statically typed object-oriented programming languages like Java and C#.

It does not apply to imperative programming languages like C where references are treated as values (e.g., with pointer arithmetic).

Scope	Toy Language Definition (Syntax)	Target Domain Component
Reference	$\text{rexp} ::= x \mid x.f \mid \text{new } C$ $\text{rcond} ::= x \diamond \text{null} \mid x \diamond y$	$\Sigma_{\text{Ref}} = \text{Env}_{\text{Ref}} \times \text{Store}_{\text{Ref}}$ (Heap Analysis, e.g., TVLA)
Value	$\text{vexp} ::= x \mid x.f \mid v_1 \diamond v_2$ $\text{vcond} ::= v_1 \diamond v_2$	$\Sigma_{\text{Val}} = \text{Env}_{\text{Val}} \times \text{Store}_{\text{Val}}$ (Value Analysis, e.g., Intervals)
Statements	$\text{st} ::= x = \text{rexp} \mid x = \text{vexp} \mid \dots$	Interaction via Substitutions

Table: Expressions, conditions and statements.

Definition (Concrete State)

The state Σ is composed of an environment and a store where types are mixed:

- **Environment:** Maps local variables to references or values.

$$\text{Env} : \text{Var} \rightarrow (\text{Ref} \cup \text{Val})$$

- **Store:** Maps heap locations (reference + field) to references or values.

$$\text{Store} : (\text{Ref} \times \text{Fld}) \rightarrow (\text{Ref} \cup \text{Val})$$

The resulting state is then the cartesian product:

$$\Sigma = \text{Env} \times \text{Store}$$

To enable generic analysis combination, we formally split the concrete domain into **two disjoint** components.

Reference State (Σ_{Ref})

Tracks only topology and pointers.

- › $\text{Env}_{\text{Ref}} : \text{Var} \rightarrow \text{Ref}$
- › $\text{Store}_{\text{Ref}} : (\text{Ref} \times \text{Fld}) \rightarrow \text{Ref}$

Value State (Σ_{Val})

Tracks only primitive values.

- › $\text{Env}_{\text{Val}} : \text{Var} \rightarrow \text{Val}$
- › $\text{Store}_{\text{Val}} : (\text{Ref} \times \text{Fld}) \rightarrow \text{Val}$

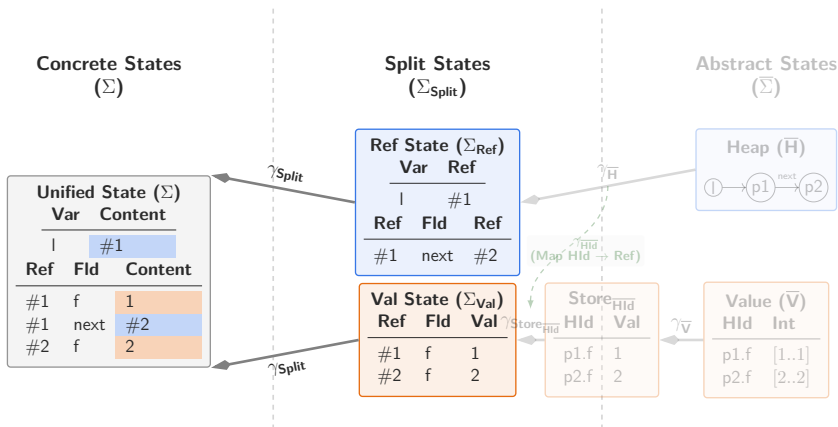
Definition (Split State Σ_{Split})

The resulting state is the cartesian product of the two components:

$$\Sigma_{\text{Split}} = \Sigma_{\text{Ref}} \times \Sigma_{\text{Val}}$$

- This domain is isomorphic to Σ but structurally separated.
- We can regard Σ_{Split} as a superstructure that is more convenient to abstract.

The Abstraction Hierarchy



Handling Mixed Expressions: Function \mathbb{R}

In the split domain, value expressions often depend on the heap (e.g., accessing $x.f$).

The Challenge

The Value Analysis (Σ_{Val}) needs to evaluate $x.f$, but it is “blind” to pointers: it does not know which reference x targets .

The Solution: Preprocessing Function \mathbb{R}

We introduce \mathbb{R} to replace local variables in field accesses with their **concrete references** from the Reference State (Σ_{Ref}).

Transformation Rules:

$$\begin{aligned}\mathbb{R}[x, (e_{\text{Ref}}, s_{\text{Ref}})] &= x \\ \mathbb{R}[x.f, (e_{\text{Ref}}, s_{\text{Ref}})] &= \langle \mathbf{e}_{\text{Ref}}(x) \rangle.f\end{aligned}$$

Integration in Semantics (Field Assignment): To evaluate $x.f = \text{vexp}$, we apply \mathbb{R} to resolve pointers on **both sides**, then delegate to Value semantics:

$$\frac{\overbrace{\langle \mathbb{R}[\![x.f]\!], \sigma_{\text{Ref}} \rangle}^{\text{Target Loc}} = \overbrace{\langle \mathbb{R}[\![\text{vexp}]\!], \sigma_{\text{Ref}} \rangle}^{\text{Resolved Expr}}, \sigma_{\text{Val}} \rangle \rightarrow_{\text{Val}} \sigma'_{\text{Val}}}{\langle x.f = \text{vexp}, (\sigma_{\text{Ref}}, \sigma_{\text{Val}}) \rangle \rightarrow_{\text{Split}} (\sigma_{\text{Ref}}, \sigma'_{\text{Val}})}$$

Lattice Structure

The Concrete Domain forms a **complete lattice** structure induced by standard set operations. The order is defined by set inclusion on the powerset of states:

$$\langle \wp(\Sigma), \subseteq \rangle$$

Like the concrete one, the lattice structure of the Split Domain is given by set of elements:

$$\langle \wp(\Sigma_{\text{Split}}), \subseteq \rangle$$

Definition (Concretization γ_{Split})

Defines how split states map back to concrete states via **pointwise set union**:

$$\gamma_{\text{Split}}(T) = \{(e_v \cup e_h, s_v \cup s_h) \mid ((e_h, s_h), (e_v, s_v)) \in T\}$$

Information

Crucial Assumption: Since the language distinguishes between value and reference expressions, the domains of e_v/e_h and s_v/s_h **do not overlap**.

Lemma (Galois Connection)

Since γ_{Split} is a complete \cap -morphism (based on set operators), it induces a valid Galois Connection:

$$\langle \wp(\Sigma), \subseteq \rangle \xLeftrightarrow[\alpha_{Split}]{\gamma_{Split}} \langle \wp(\Sigma_{Split}), \subseteq \rangle$$

where $\alpha_{Split} = \lambda X. \cap \{Y : X \subseteq \gamma_{Split}(Y)\}$.

Proof.

α_{Split} is well-defined since γ_{Split} is a complete \cap -morphism since it is based on set operators. The fact that it forms a Galois connection follows immediately from the definition of α_{Split} [CC77]. □

1 Introduction

- The problem

2 The framework

- Concrete domain and semantics
- Abstract domain and semantics
- Instantiation

3 Conclusion and Contributions

The framework is designed to be **parametric**. We assume the existence of two abstract domains interacting to approximate the Split State.

1. Heap Analysis (\bar{H})

Approximates the reference state Σ_{Ref} .

- › Lattice: $\langle \bar{H}, \sqsubseteq_{\bar{H}}, \sqcup_{\bar{H}}, \sqcap_{\bar{H}} \rangle$
- › Widening: $\nabla_{\bar{H}}$

2. Value Analysis (\bar{V})

Approximates the value state Σ_{Val} .

- › Lattice: $\langle \bar{V}, \sqsubseteq_{\bar{V}}, \sqcup_{\bar{V}}, \sqcap_{\bar{V}} \rangle$
- › Widening: $\nabla_{\bar{V}}$

Definition (Abstract State $\bar{\Sigma}$)

The combined abstract state is the Cartesian product:

$$\bar{\Sigma} = \bar{H} \times \bar{V}$$

equipped with pointwise lattice operators.

Connecting Heap and Value: Heap Identifiers

The two domains need a common language to communicate.

- **Heap Domain (\overline{H}):** Knows about pointers, structure, and reachability. It abstracts memory addresses (\mathbb{A}).
- **Value Domain (\overline{V}):** Knows about numerical properties (intervals, polyhedra), but operates on a set of *variables*, not addresses.

Connecting Heap and Value: Heap Identifiers

The two domains need a common language to communicate.

- **Heap Domain (\overline{H}):** Knows about pointers, structure, and reachability. It abstracts memory addresses (\mathbb{A}).
- **Value Domain (\overline{V}):** Knows about numerical properties (intervals, polyhedra), but operates on a set of *variables*, not addresses.

The Solution: Abstract Heap Identifiers

The Heap Domain exports a set of symbolic names, called **Heap Identifiers**.

$$\overline{\text{HId}} = \{\text{id}_1, \text{id}_2, \dots, \text{id}_n\}$$

- To \overline{H} , they represent abstract nodes (or regions of memory).
- To \overline{V} , they are treated simply as **variables**.

Since the Value Analysis operates on identifiers, its concretization produces stores with abstract heap identifiers instead of concrete locations.

Abstract Store on Identifiers

$$\text{Store}_{\overline{\text{HId}}} : \overline{\text{HId}} \rightarrow \wp(\text{Val})$$

Codomain is a set of values because a summary node maps to many concrete references with different values.

Definition ($\gamma_{\overline{\text{V}}}$)

The concretization of an abstract value state $\bar{v} \in \overline{\text{V}}$ produces:

- 1 Environments in Env_{Val}
- 2 Stores in $\text{Store}_{\overline{\text{HId}}}$

Missing link: We have values for identifiers, but we need to know *where* these identifiers are in memory.

The Heap Analysis tracks the shape and symbolically represents nodes using $\overline{\text{HId}}$. To concretize, it must map these symbols to concrete locations.

The Mapping Function $\gamma_{\overline{\text{HId}}}$

The heap concretization provides a function relating identifiers to concrete locations ($\text{Ref} \times \text{Fld}$):

$$\gamma_{\overline{\text{HId}}} : \overline{\text{HId}} \rightarrow \wp(\text{Ref} \times \text{Fld})$$

Definition ($\gamma_{\overline{\text{H}}}$)

Formally, $\gamma_{\overline{\text{H}}} : \overline{\text{H}} \rightarrow \wp(\Sigma_{\text{Ref}} \times (\overline{\text{HId}} \rightarrow \wp(\text{Ref} \times \text{Fld})))$

It returns a set of pairs containing:

- A concrete state σ_{H} .
- A concretization of heap identifiers $\gamma_{\overline{\text{HId}}}$ to compute $\gamma_{\overline{\text{V}}}$.

Finally, we define the concretization of the abstract state (\bar{h}, \bar{v}) into the Split domain.

Definition ($\gamma_{\overline{\Sigma}}$)

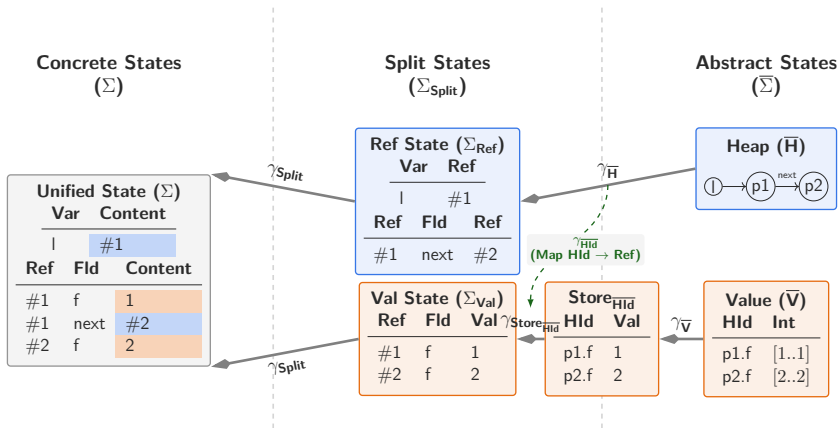
Formally, $\gamma_{\overline{\Sigma}}((\bar{h}, \bar{v})) : \overline{\Sigma} \rightarrow \wp(\Sigma_{\text{Split}})$

$$\gamma_{\overline{\Sigma}}((\bar{h}, \bar{v})) = \{(\sigma_H, (e_v, s_v)) \mid \text{conditions hold}\}$$

where:

- ▶ $(e_v, s_{\overline{\text{HId}}}) \in \gamma_{\overline{V}}(\bar{v})$
(Value Analysis provides Env_{Val} and $\text{Store}_{\overline{\text{HId}}}$)
- ▶ $(\sigma_H, \gamma_{\overline{\text{HId}}}(\bar{h})) \in \gamma_{\overline{H}}(\bar{h})$
(Heap Analysis provides a state in Σ_{Ref} and ID-map)
- ▶ $s_v \in \gamma_{\text{Store}_{\overline{\text{HId}}}}(s_{\overline{\text{HId}}}, \gamma_{\overline{\text{HId}}}(\bar{h}))$
(Combined to get concrete $\text{Store}_{\text{Val}}$)

The Abstraction Hierarchy



Remark (Requirements for Soundness (Conditions C1–C5))

The validity of the following Lemma relies on strict conditions on

γ_{HId} :

- **Structural Consistency (C1, C3, C5):**
 - All heap identifiers must be concretizable (C1).
 - Identifiers must represent *disjoint* memory portions to allow independent updates (C3).
 - Non-summary nodes must map to exactly one reference (C5).
- **Mathematical Properties (C2, C4):** These ensure γ_{Σ} is **meet-preserving**. The heap identifiers' concretization of the intersection of heaps is the pointwise intersection of the heap identifiers' concretization of all the intersected states.

Lemma (Galois Connection)

Since $\gamma_{\bar{\Sigma}}$ is a complete \cap -morphism (based on set operators), it induces a valid Galois Connection:

$$\langle \emptyset(\Sigma_{Split}), \subseteq \rangle \xrightleftharpoons[\alpha_{\bar{\Sigma}}]{\gamma_{\bar{\Sigma}}} \langle \bar{\Sigma}, \sqsubseteq \rangle$$

where $\alpha_{\bar{\Sigma}} = \lambda X. \cap \{Y : X \subseteq \gamma_{\bar{\Sigma}}(Y)\}$.

In the abstract domain, the Value Analysis (\overline{V}) tracks information on **Heap Identifiers** (abstract nodes), not on raw pointer paths.

The Abstract Solution

We define the abstract preprocessing function $\overline{\mathbb{R}}$ to translate field accesses into the corresponding **set of Heap Identifiers** provided by the Heap Analysis (\overline{H}).

Unlike the concrete case, this translation may return **multiple** results (e.g., if a variable points to multiple abstract nodes).

Abstract Transformation Rules :

$$\begin{aligned}\overline{\mathbb{R}}[x, \overline{h}] &= \{x\} \\ \overline{\mathbb{R}}[x.f, \overline{h}] &= I \quad \text{where } \langle x.f, \overline{h} \rangle \rightarrow_{\overline{H}} I\end{aligned}$$

Note: I is the set of heap identifiers (e.g., $\{id_1, id_2\}$) retrieved by the Heap Analysis.

Semantics Integration (Updates)

Once $\overline{\mathbb{R}}$ resolves the field access, the Value Analysis performs the assignment. Formally, we compute the join of all possible outcomes for each identifier i returned by $\overline{\mathbb{R}}$:

$$\overline{v}_{new} = \bigsqcup \{ \overline{v}' \mid i \in \overline{\mathbb{R}}[[x.f, \overline{h}]], \langle i = \dots, \overline{v} \rangle \rightarrow_{\overline{V}} \overline{v}' \}$$

The final state depends on the precision of the Heap Analysis:

- **Strong Update:** If i is unique and definite ($\neg \text{isSum}(i)$), the old value is replaced:

$$\overline{v}_{post} = \overline{v}_{new}$$

- **Weak Update:** If i is a summary node or multiple identifiers exist, we must join with the previous state to preserve soundness:

$$\overline{v}_{post} = \overline{v} \sqcup_{\overline{V}} \overline{v}_{new}$$

1 Introduction

- The problem

2 The framework

- Concrete domain and semantics
- Abstract domain and semantics
- Instantiation

3 Conclusion and Contributions

The framework is validated by plugging in two different heap abstractions:

1 **Pointer Analysis (PA)**

A flow-sensitive analysis based on *Allocation Site Abstraction* [And94; MSV10].

» **Heap Identifiers:** Defined by the allocation site label and the field name:

$$\overline{\text{HId}}_{\text{PA}} = \text{Lab} \times \text{Field}$$

» **Substitutions:** Since abstract locations are statically determined, no dynamic splitting occurs. Therefore, substitutions are always empty (\emptyset).

2 TVAL+ (Shape Analysis)

Based on the 3-valued logic engine TVLA [SRW02].

- » **Identity:** Introduces *Name Predicates* to track nodes across transformations.
- » **Normalization:** A merge/split procedure ensures states are normalized.
- » **Substitutions:** Generated dynamically to inform the value analysis about node materialization and summarization.

The generic framework allows standard numerical domains to track information over heap contents.

3 Numerical Domains

Standard domains (e.g., Intervals, Octagons) are plugged into the value component.

- » **Transparency:** Heap Identifiers ($\overline{\text{HId}}$) are treated transparently, exactly as local variables.
- » **Handling Summaries:** To preserve soundness, assignments to summary nodes ($\overline{\text{isSum}}(id) = \text{true}$) are handled via **weak updates** (join of old and new values).

Conclusion on Instantiation: The framework successfully bridges static approaches (PA) and dynamic approaches (TVLA) with numerical reasoning without losing soundness.

1 Introduction

- The problem

2 The framework

- Concrete domain and semantics
- Abstract domain and semantics
- Instantiation

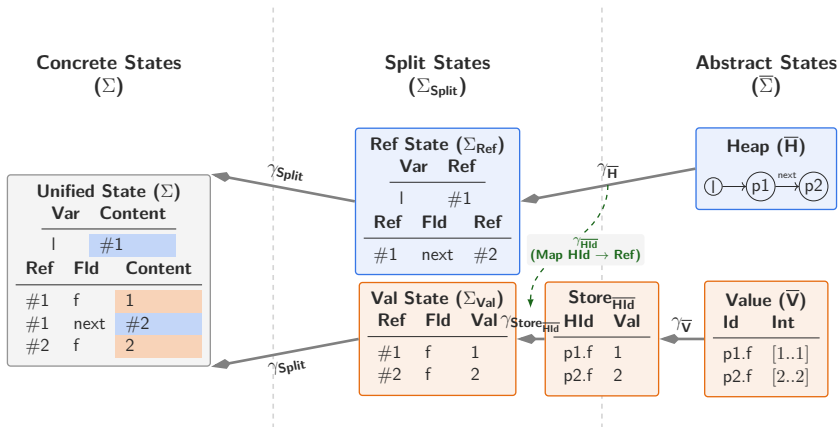
3 Conclusion and Contributions



The paper presented a formal approach to the static analysis of object-oriented languages.

Key Contributions:

- **Generic Framework:** A unified formalization to automatically combine arbitrary Heap and Value analyses.
- **Handling Dynamics:** The *first* generic framework capable of supporting **materialization** and **summarization** of heap identifiers.
- **Soundness:** Proved soundness of the combination, relying on standard Abstract Interpretation operators and a specific interface (Substitutions).
- **Versatility:** Successfully instantiated with:
 - Standard Pointer Analysis (Static).
 - TVAL+ / TVLA (Dynamic Shape Analysis).
 - Numerical Domains (Value Analysis).



Thank you for the attention!

- [And94] Lars Ole Andersen. "Program analysis and specialization for the C programming language". PhD thesis. DIKU, University of Copenhagen, 1994.
- [CC77] Patrick Cousot and Radhia Cousot. "Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints". In: *Proceedings of the 4th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. POPL '77. ACM, 1977, pp. 238–252.
- [MSV10] Matthew Might, Yannis Smaragdakis, and David Van Horn. "Resolving and exploiting the k-CFA paradox". In: *Proceedings of PLDI '10*. ACM. 2010.
- [SRW02] Mooly Sagiv, Thomas Reps, and Reinhard Wilhelm. "Parametric shape analysis via 3-valued logic". In: *ACM Trans. Program. Lang. Syst.* 24.3 (2002), pp. 217–298.