

LEARNING TO REPRESENT AND APPLY EDITS

Background Edits are ubiquitous in both source code and language

Before: `myDict[key] = obj.Value` Before: we **have** a great day
 After: `myDict.Add(key, obj.Value)` After: we **had** a great day
 source code edits (e.g., GitHub commits) natural language edits

Research Questions

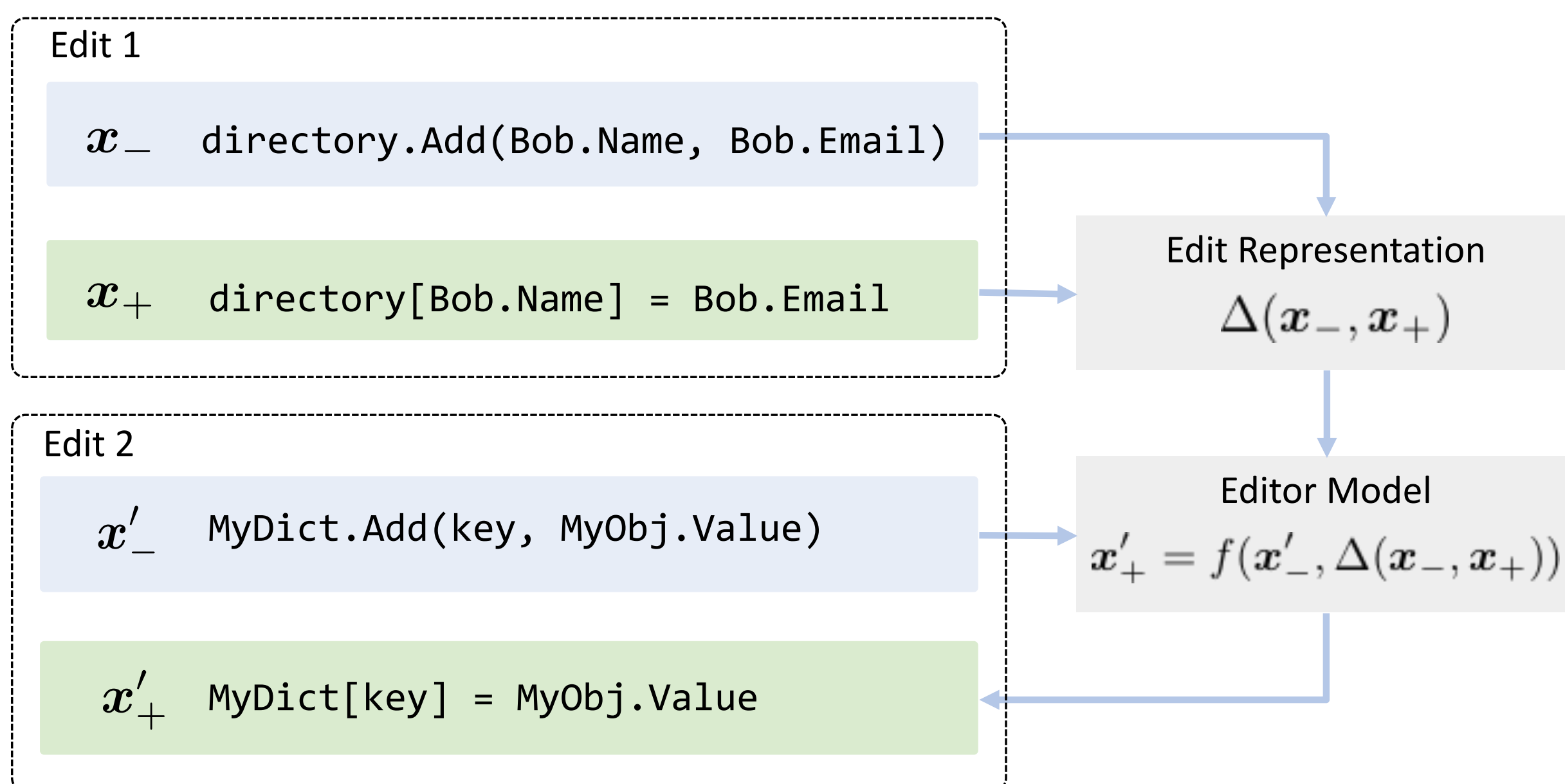
Edit Representation: How to learn edit patterns from data?

Neural Editor: How to apply an edit to a new input?

Contributions

- Formalize the problem of data-driven learning of edits.
- Present neural models that learn to represent and apply edits.
- Release a large-scale dataset of code edits for future research.

Key Idea – End-to-End Learning to Edit



Learning Objective maximize the likelihood of generating the target x'_+ given the input x'_- and the “ground-truth” edit $\Delta(x'_-, x'_+)$

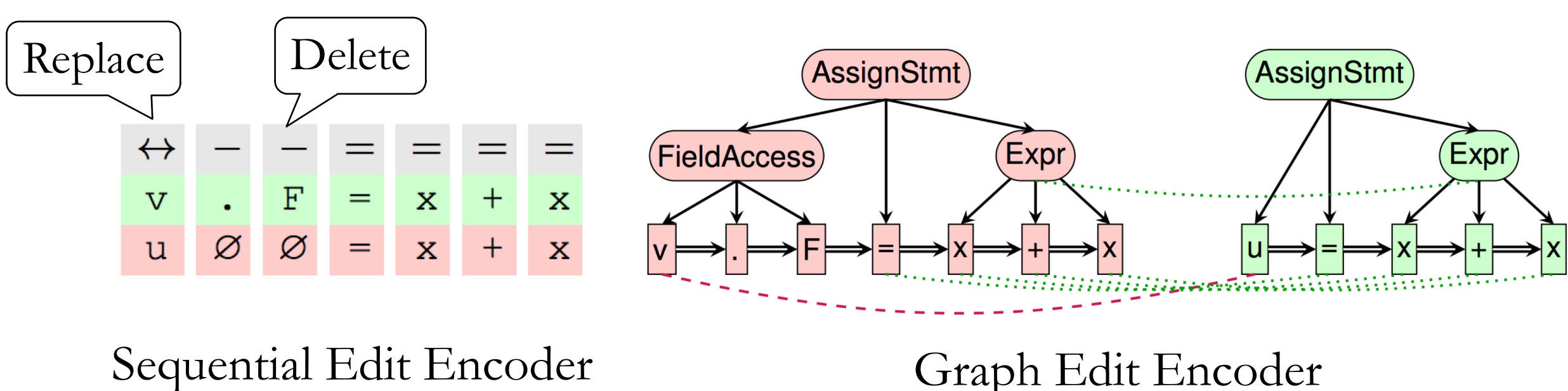
EDIT REPRESENTATIONS

Sequential Encoding of Edits

- Use deterministic diffing algorithm to get alignments of tokens in the source and target, which are encoded using bi-LSTMs

Graph Encoding of Edits

- Source and target data is represented as syntax trees
- Add alignment edges between edited nodes. Encode the graph using graph neural networks (Li et al., 2016)

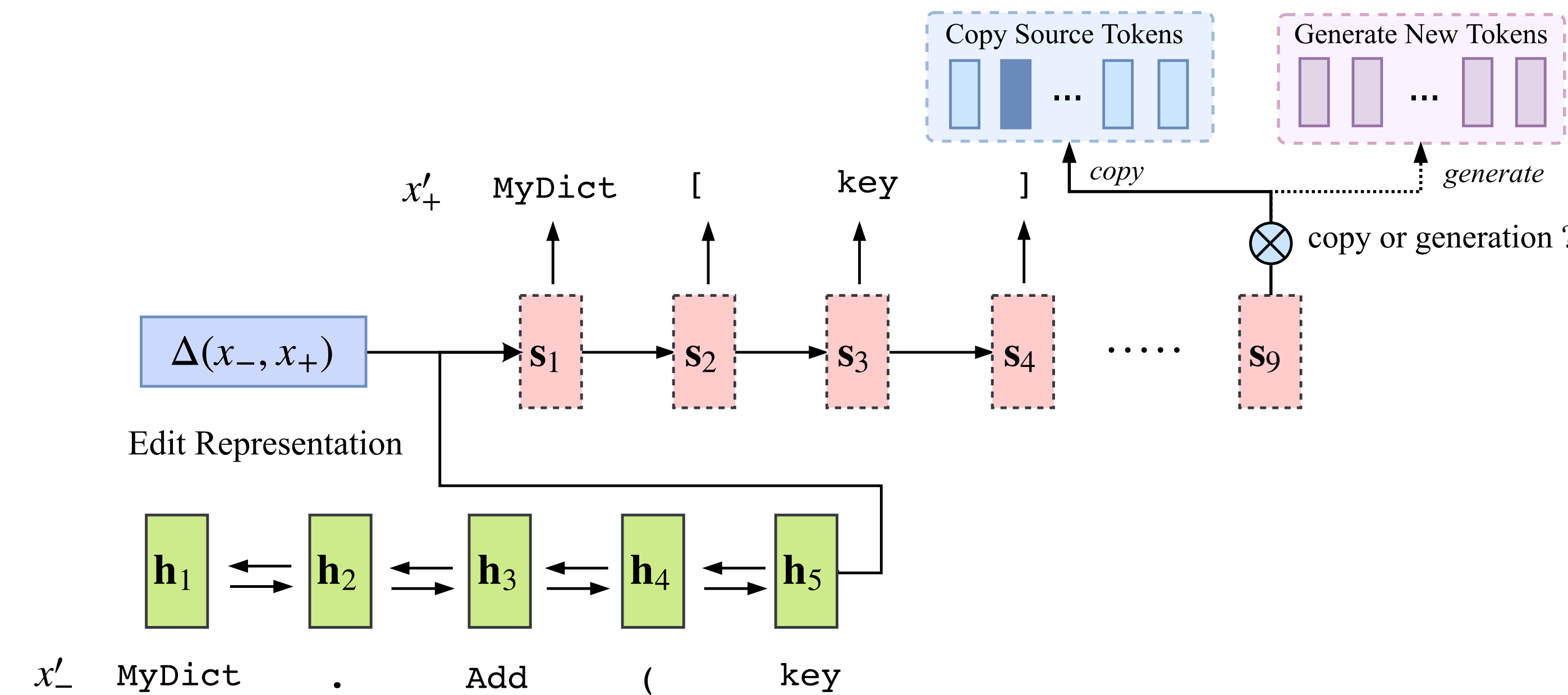


NEURAL EDITORS

Overview Given input x'_- , and edit representation $\Delta(x_-, x_+)$, a neural editor applies the edit to x'_- and generates the updated input x'_+

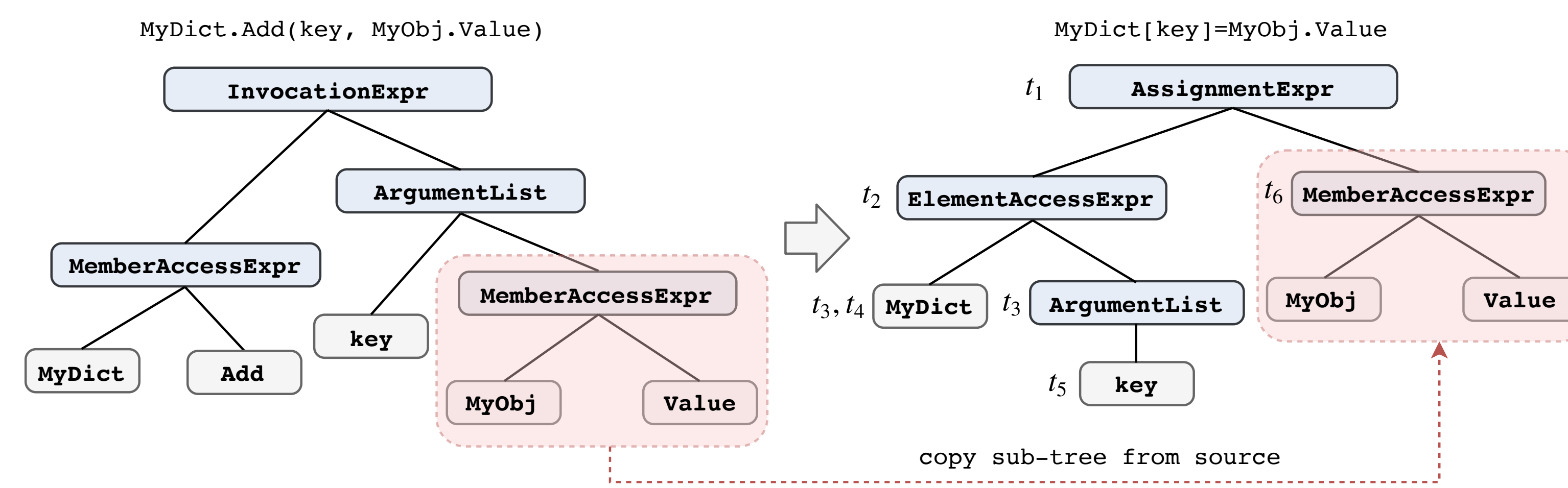
Sequence-to-Sequence Neural Editors

- Seq2Seq editors encode the input data as a sequence of tokens
- A recurrent neural network generates the target x'_+ using the input x'_- and the edit representation $\Delta(x_-, x_+)$



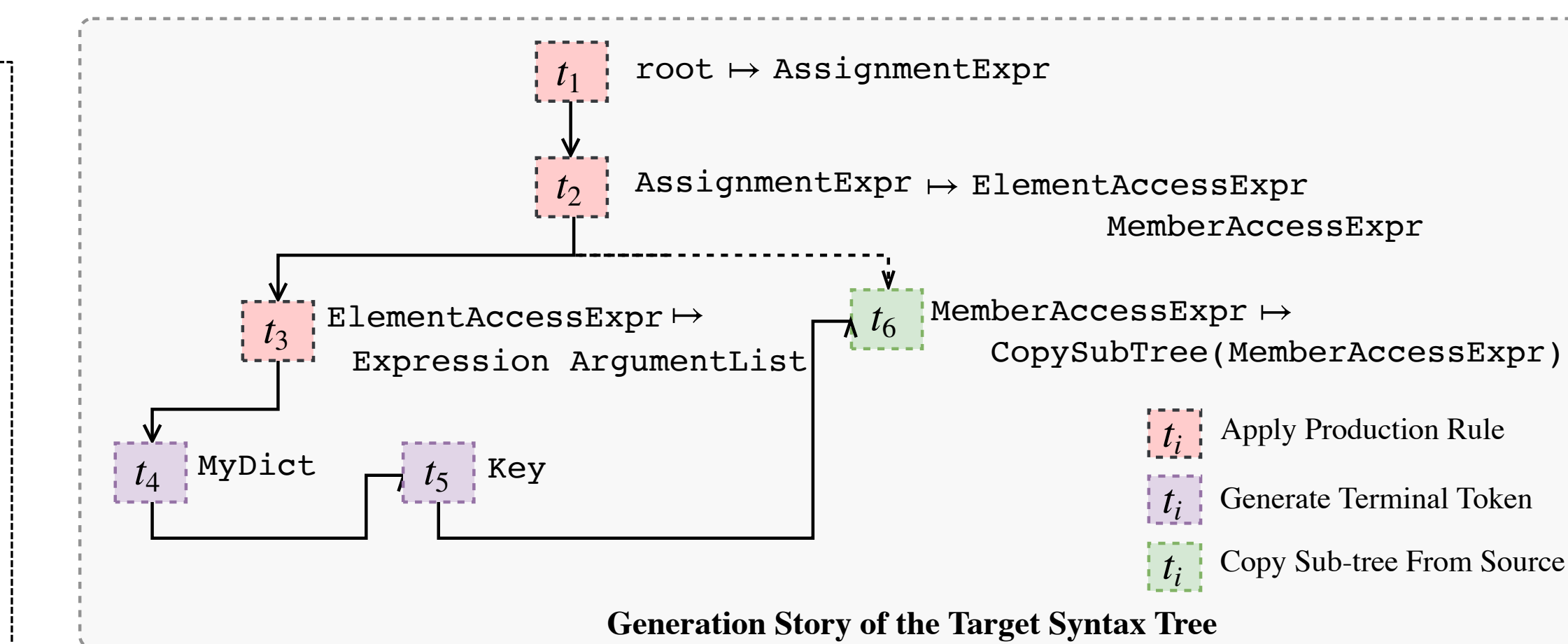
Graph-to-Tree Neural Editors

Motivation the data to edit (e.g., source code) usually has strong underlying structure. How to utilize the structural information to better predict the edited output?



Graph2Tree editors transduce the tree-structured input to the target tree via a sequence of tree-constructing actions (Yin et al., 2018)

- Source and target data (code) is represented as structured abstract syntax trees
- Source data is encoded using gated graph neural networks (Allamanis et al., 2018)



GITHEEDIT CODE EDIT DATASET

- 110K pairs of C# code edits collected from GitHub commit histories
- Each edit involves at most three consecutive lines of code
- Shipped with parsed abstract syntax trees for each sample



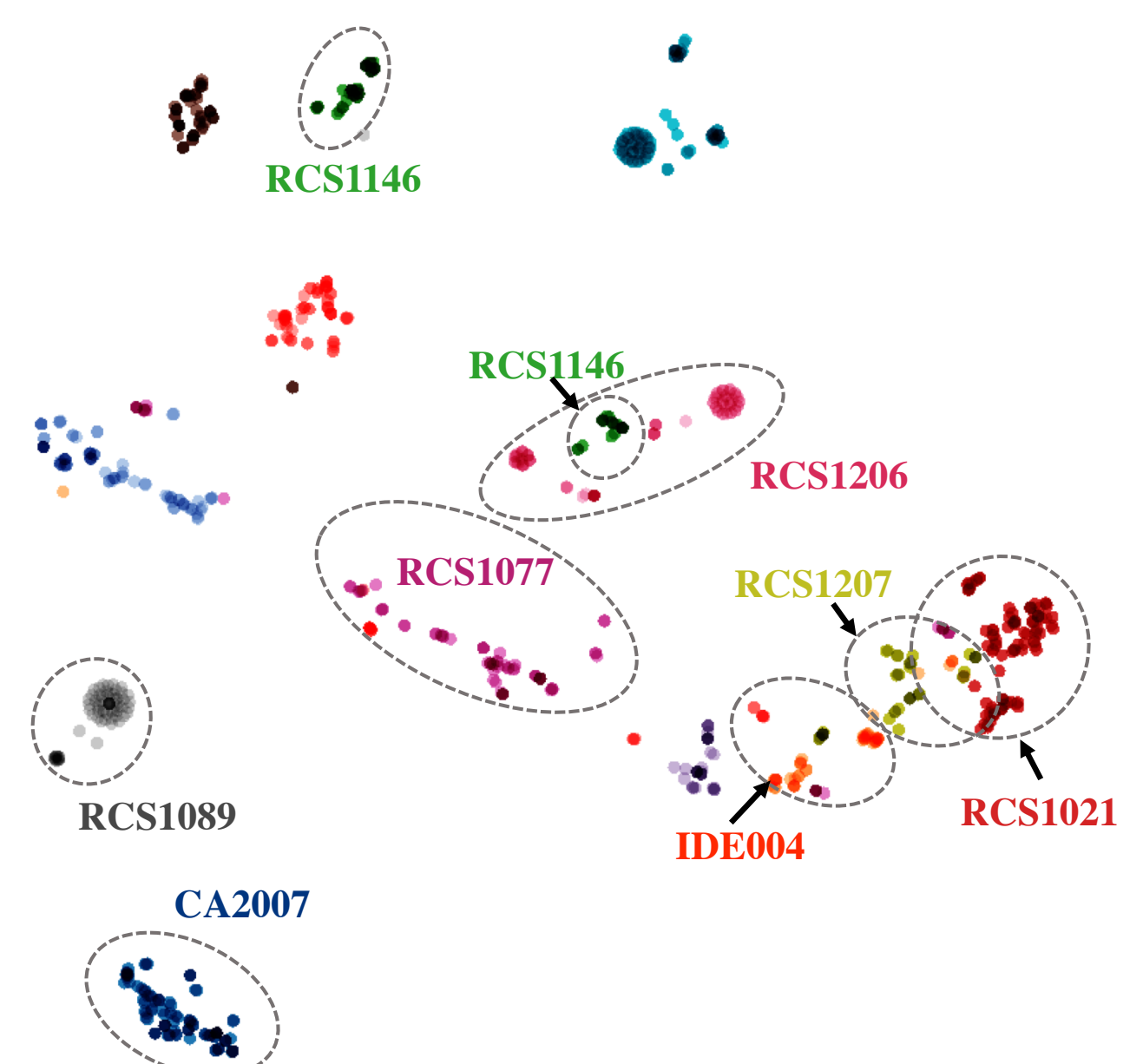
EXPERIMENTS

Quality of Edit Representations

1. t-SNE Visualization Systems trained on the noisy GITHUBEDITS, and tested on 2.8K edit pairs with known, labeled edit categories (16 in total). Examples:

Before `x != null && x.StartsWith("a")`
 After `x?.StartsWith("a")`
 Category Use conditional Access (RCS1146)

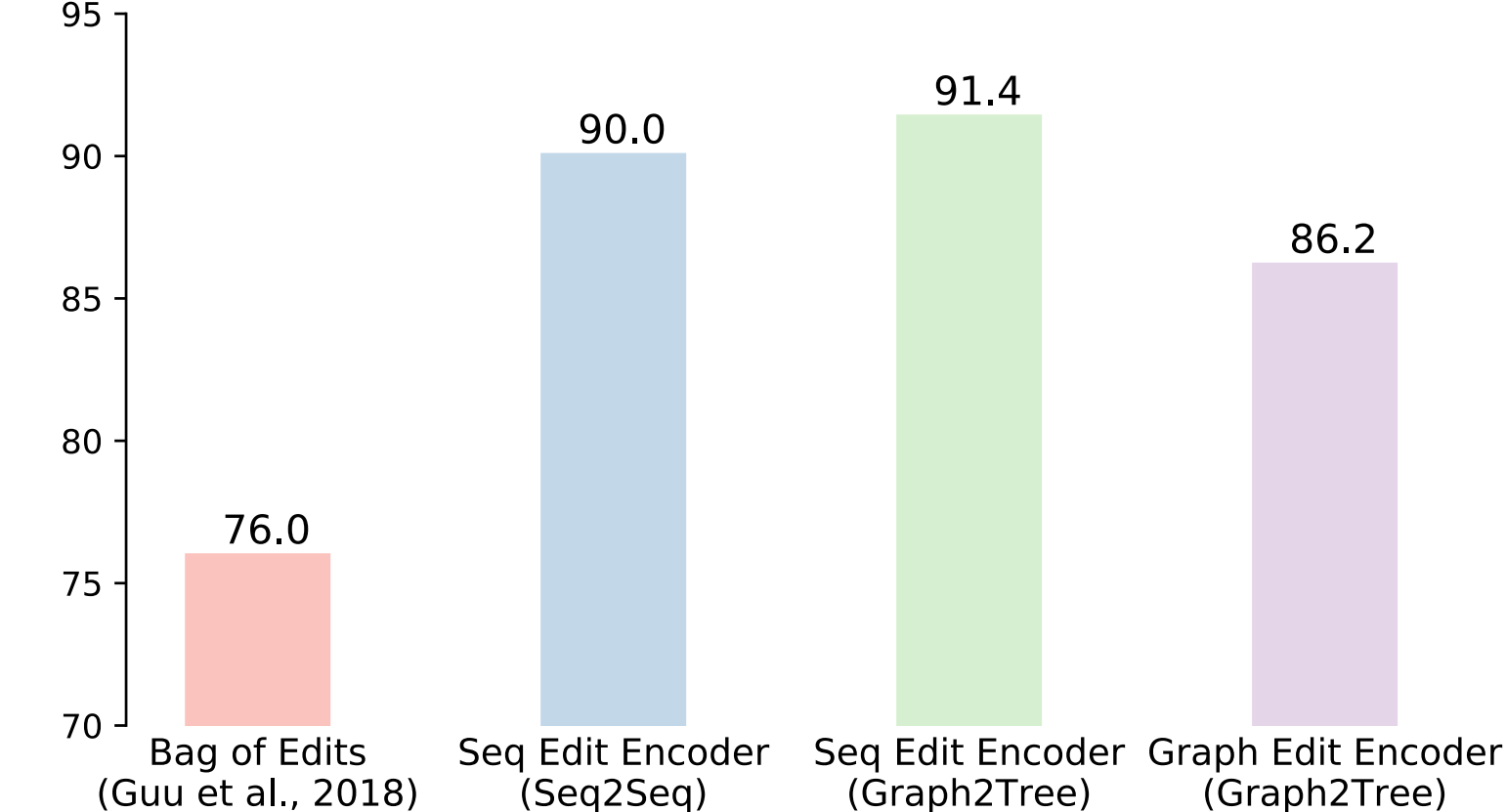
Before `i = i + 1`
 After `i++`
 Category Use --/++ operator (RCS1089)



2. Edit Retrieval Task Given a seed edit representation, retrieve and evaluate the relevance of retrieved neighboring edits

Manually annotate the relevance with a scale of 3 (*highly relevant, relevant, irrelevant*)

NDCG@3 on Edit Retrieval Task

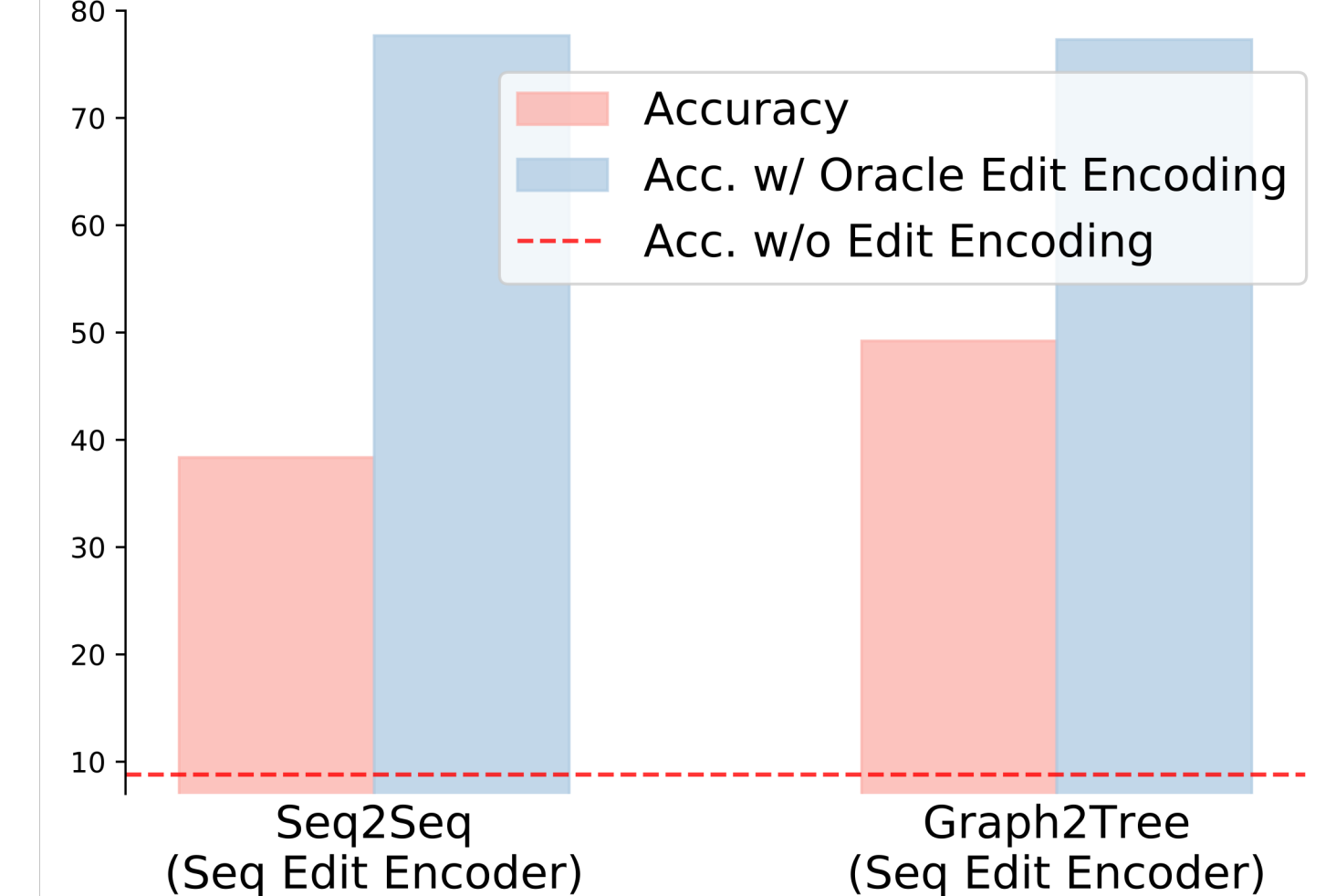


Precision of Neural Editors

End-to-End Evaluation Given a learned edit encoding $\Delta(x_-, x_+)$, apply the edit to a similar input x'_- to generate the edited input x'_+

Compare with the upper-bound accuracy of using the oracle edit encoding $\Delta(x'_-, x'_+)$

Accuracy of Neural Editors



Clustering Edits on GitHub Commits and Wikipedia Edit History

<p>Desc. Optimize LINQ Queries</p> <p>Before <code>V0.Customers.Where(V1 => V1.CustomerID == LITERAL).FirstOrDefault()</code></p> <p>After <code>V0.Customers.FirstOrDefault(V1 => V1.CustomerID == LITERAL)</code></p>	<p>Desc. Switch from Assert.Equal to Assert.Empty</p> <p>Before <code>Assert.Equal(0, V0.ProjectIds.Count)</code></p> <p>After <code>Assert.Empty(V0.ProjectIds)</code></p>
<p>Before <code>this.V0.Where(V1 => V1.CanDeserialize(V2)).FirstOrDefault()</code></p> <p>After <code>this.V0.FirstOrDefault(V1 => V1.CanDeserialize(V2))</code></p>	<p>Before <code>Assert.Equal(0, V0.ProjectReferences.Count())</code></p> <p>After <code>Assert.Empty(V0.ProjectReferences)</code></p>
<p>Before <code>V0.TypeConverters.Where(V1 => V1.CanConvertTo(V2, V0)).FirstOrDefault()</code></p> <p>After <code>V0.TypeConverters.FirstOrDefault(V1 => V1.CanConvertTo(V2, V0))</code></p>	<p>Before <code>Assert.Equal(0, V0.Messages.Count)</code></p> <p>After <code>Assert.Empty(V0.Messages)</code></p>

Sampled edit clusters from GITHUBEDITS (omitting variable assignments for clarity)

<p>Desc. Add a parenthetical expression <i>also ... as mid-state regional airport, +also known as mid-state airport,+ is a small airport on in rush township...</i></p> <p>Before <code>islaamic culture, +also known as saracenic culture,+ is a term primarily used in secular academia...</code></p> <p>After <code>birds of prey, +also known as raptors,+ are birds that hunt for food primarily via flight...</code></p>	<p>Desc. Add a Person's Middle Name</p> <p>Before <code>isaiah +marcus+ rankin (born 22 may 1978 in london) is an english professional footballer</code></p> <p>After <code>audrey +Kathleen+ brown (born 24 may , 1913) is a british athlete who competed mainly in the 100 metres .</code></p> <p>Before <code>monique +edith+ lamoureux (born july 3 , 1989) is an american ice hockey player .</code></p>
---	---

Sampled edit clusters (+insertion+ edits) from WIKIATOMICEDITS dataset (Faruqui et al., 2018)