² SNAP, Stanford University

tl;dr

- We set the new SOTA on code summarization by learning from Structure and Context of programs.
- Our language-agnostic design enables joint training on multiple programming languages.
- Our novel **multi-language training** substantially improves performance on each individual language.

Task: Code Summarization

- Predict a method's name given its body.
- A common task to evaluate ML models' performance on code.
- In high-quality code bases, a method's name summarizes its functionality.
- □ Labels come 'for free', and there is lots of high-quality open-source code.

Example



Code Transformer



ICLR 2021

Language-Agnostic Representation Learning of Source Code from Structure and Context

Daniel Zügner¹, Tobias Kirschstein¹, Michele Catasta², Jure Leskovec², Stephan Günnemann¹

(0.8) parse_bool to_lower

_(0.01) get_string

Structure and Context Representation

Structure and Context are complementary representations of a program.

Source code as sequence of tokens (referred to as Context)





Typical LSTMs, Transformers models:

GNNs, LSTMs on AST paths

- □ For instance, get_model() is far away from the return statement in the Context; it is hard for a model to capture such long-range dependencies.
- □ In the AST, the two corresponding nodes are only a few hops apart.
- Many previous works only learn from either Structure or Context.
- □ We propose the Code Transformer, which learns jointly from both.

Relative Distances on the AST

- Can be computed from the AST for any programming language. • Capture the local and global structure in the AST.
- **Language-agnostic design** instead of specialized proprietary pipelines.
- Enables training on any programming language, even jointly (see results).
- Efficient binning enables using continuous-valued distances (e.g., PPR).

May 3-7 Virtual Event

Source code as Abstract Syntax Tree (AST; referred to as Structure)

Results

	Model	Python	Javascript	Ruby	Go
e- age	code2seq	29.3	24.0	14.3	47.5
eng gug	GREAT	33.2	28.9	23.4	48.2
Si lan	Code Transformer	35.0	32.1	27.5	51.3
e.	code2seq*	29.3	26.1	19.9	48.2
nag	Change (Δ)	-0.02	+2.09	+5.65	+0.72
ngu	GREAT*	33.9	30.4	26.9	50.4
	Change (Δ)	+0.70	+1.55	+3.51	+2.18
Jult	Code Transformer	36.2	33.2	31.2	53.0
2	Change (Δ)	+1.21	+1.08	+3.74	+2.63
Mult. +	Code Transformer	37.4	34.3	32.0	54.7
Pretrain	Change (Δ)	+2.38	+2.20	+4.51	+3.31

* Adapted by us for multilanguage training

Strongest gains are on Ruby, the language with fewest training samples.

Ablation study

Table: Ablation of Structure and Context (Java-small dataset)

ST distance	F1 score	AST distance
ode Transformer + Pretrain	53.77	Sibling shortest paths
Code Transformer	52.22	Ancestor shortest paths
– Structure	50.34	Shortest paths
– Context	49.45	Personalized PageRank
 Pointer Network 	48.50	All of the above

Multilanguage embeddings

- We obtain a shared embedding space for multiple programming languages.
- Similar methods tend to be mapped to regions close-by in embedding space.
- We can 'measure' similarity of programming languages via their embeddings.

Figure: t-SNE visualization of multilanguage embeddings

Demo: http://code-transformer.org

Stanfor University



Table: Code Summarization Results on the CodeSearchNet dataset (F1 score).

• Our Code Transformer outperforms all baselines in all scenarios.

Multilanguage training improves results on each individual language.

□ This is true for all studied models using the Structure, opening exciting directions for future research. No improvement for Context-only models.

Table: Ablation of AST distances w/o Context (Java-small dataset)



Code: www.daml.in.tum.de/code-transformer