

Robust representation of small organic molecules with self-referencing embedded strings (SELFIES)

Florian Häse

Data Scientist – Computational Life Science

Bayer AG, Division Crop Science





Describing a chemical substance

Which representation is most useful?



This chemical with this structure!



Different descriptions of the same chemical substance might be helpful in different situations.

This DB entry with these identifiers!

Representations for molecular modeling

Which aspects of the molecular are most relevant to a given modeling task?



3

// Many different representations have been developed for molecular modeling tasks

Each representation emphasizes certain aspects of a molecule, and obfuscates other aspects

- # Every modeling task should start with the questions
 - which aspects about the molecule are fundamental to the task?
 - // which representation highlights these aspects best?

Example: Coulomb matrices to predict electronic properties

The **Coulomb Matrix** is a *global descriptor* based on the *pairwise electrostatic interaction* between nuclei,

$$M_{ij}^{\text{Coulomb}} = \begin{cases} Z_i^{2.4}/2 & \text{for } i = j \\ Z_i Z_j / R_{ij} & \text{for } i \neq j \end{cases}$$

where

 Z_i ... nuclear charge R_{ij} ... interatomic distance

which

- // uniquely encodes molecule structures
- // satisfies translational and rotational invariance



Task:

Predict atomization/formation energies for QM7b



Multidimensional distributions of interatomic many-body expansions (Faber *et al.*) yield much more accurate prediction results

Better representations produce better results

Faber, Felix A., et al. "Alchemical and structural distribution-based representation for universal quantum machine learning." The Journal of chemical physics 148.24 (2018): 241717.

Rupp, Matthias, et al. "Fast and accurate modeling of molecular atomization energies with machine learning." Physical review letters 108.5 (2012): 058301.



Molecules as text

Sequential approaches to molecular structures



De-novo generation of small organic molecules

How can we generate molecule structures from scratch?

Goals:

- // machine-generated molecule structures
- // which satisfy basic chemical constraints

"valid" molecule



"invalid" molecules

Idea:

- get inspired by educational toys *chemistry sets*!
- // accept the constraints of your scientific kit
- // only attach those atoms that fit the current structure

Example:



Algorithm (ignoring H):

- // Take one carbon (C)
- // Take another carbon (C)
- // Take one oxygen (O)
 - $C \rightarrow CC \rightarrow CCO$

Molecule can be "written as a string"

Formalizing the generation of molecular structures

Which set of rules enables the sequential generation of molecular graphs?

Generation process is governed by rules and representations



Formal grammars

are composed of

- // terminal symbols T
- // nonterminal symbols **N**
- // production rules P
- // start symbols S

7

ExampleProduction rules: $T = \{\mathbf{C}, \mathbf{O}\}$ // (i) [molecule] $\rightarrow >>$ empty <<</td> $N = \{[molecule]\}$ // (ii) [molecule] $\rightarrow \mathbf{C}$ [molecule] $S = \{[molecule]\}$ // (iii) [molecule] $\rightarrow \mathbf{O}$ [molecule]

Execution of, e.g., [(ii), (ii), (ii), (i)] leads to [molecule] -> C [molecule] -> C C [molecule] -> C C O [molecule] -> C C O

Observations:

- // This grammar generates, e.g., alkanes and alcohols CCO, CCCC, OCCC, C
- // This grammar cannot generate all molecules CNO is a syntactical mistake (N is not part of the vocabulary)
- Some of the generated strings are not molecules OOOO is not a valid molecule



Short interlude:

Simplified Molecular Input Line Entry System (SMILES)



Excursion to SMILES strings

What is the *de-facto* standard for text-based molecular representations?

SMILES:

9

Simplified Molecular Input Line Entry System

Atoms: C, O, F, N, P, S, ... Bonds: =, #, :, -Branches: (,) Rings: 1, 2, 3, ...

Example: 3,4-methylenedioxymethamphetamine (MDMA)



Benefits

- capable of representing molecules as strings
- molecular representations are intuitive

Disadvantage

arbitrary sequences of SMILES characters do not necessarily encode valid molecules

Examples: Syntactic violations

CC1CO1CCC2CCCRing not closedCC(CCO(CCC)CCBranch not closed

Examples: Semantic violations

C=C**=O=**CCOC C=CCC**F**OC

Incorrect valency Incorrect valency



Self-referencing embedded strings (SELFIES)

Requirements on grammars for sequential molecular generation

Which properties does a grammar need to satisfy to enable sequential molecular generation?

Why did we fail **syntactically**?

CC1CO1CCC2CCC CC(CCO(CCC)CC Ring not closed Branch not closed

Some generated symbols cannot be interpreted in the given context

Mitigation: Grammar must be context-free

Why did we fail **semantically**?

C=C=O=CCOC C=CCCFOC Incorrect valency Incorrect valency

Some generated symbols do not satisfy physical / chemical constraints

Mitigation: Grammar must account for *chemistry*

Categorizing formal grammars (Chomsky)

11

Class	Grammars	Automaton
Type 0	unrestricted	Turing machine
Type 1	context sensitive	Linear-Bound
Type 2	context free	Pushdown
Туре 3	regular	Finite

Chomsky Type 2 Grammar

- // uses non-terminal symbols Z and terminal symbols t
- ${}^{\prime\prime}$ uses production rules of the form Z
 ightarrow t ~Z
- // can apply production rules regardless of the context
- but does *a priori* not produce chemically relevant sequences

Simple Chomsky type-2 grammar to generate FC=C=N

How do we design Chomsky type-2 grammars for molecular generation?

Simple grammar				[F] $[C]$ $[N]$
Terminals $T = \{ F, C, =N, =C \}$ Non-terminals $N = \{Z_0, Z_1, Z_2, Z_3\}$ Start symbol $S = \{Z_0\}$	Production rules $ \begin{array}{cccc} & & [F]: & Z_0 \rightarrow \mathbf{F} & Z_1 \\ & & [C]: & Z_1 \rightarrow \mathbf{C} & Z_3 \\ & & [N]: & Z_2 \rightarrow = \mathbf{N} \\ & & & [C]: & Z_3 \rightarrow = \mathbf{C} & Z_2 \end{array} $	Pushdown automaton generates string based on collection of specific production rules	State of derivation S_1 Z_1 Z_2 Z_3	FZ_{1} CZ_{3} $=N$ $=CZ_{2}$

Example: Consider the collection of production rules

[F] [C] [C] [N]

Starting from Z_0 , the automaton generates:

$$Z_0 \xrightarrow{[F]} \mathbf{F} Z_1 \xrightarrow{[C]} \mathbf{F} \mathbf{C} Z_3 \xrightarrow{[C]} \mathbf{F} \mathbf{C} = \mathbf{C} Z_2 \xrightarrow{[N]} \mathbf{F} \mathbf{C} = \mathbf{C} = \mathbf{N}$$

Keeping track of the current state, the pushdown automaton will only generate chemically valid structures if the production rules satisfy these constraints

Introducing SELFIES

Which features do we need to generate molecular graphs?

We require a set of decoding rules to reliably write (molecular) graphs from character-based sequences

- // writing **vertices:** non-terminal Z_r and rule A_n **generates** Vertex $t_{r,n}$ and non-terminal $Z_{h_{r,n}}$
- // writing **rings:** non-terminal Z_r and rule A_{n+m} closes N-membered ring and non-terminal $Z_{i_{r,n}}$
- // writing **branches:** non-terminal Z_r and rule A_{n+n+p} opens N-membered branch and non-terminal $Z_{h_{r,n}}$

However: These decoding rules do not account for *chemical / physical semantics* (yet)!

Introducing SELFIES

Which features do we need to generate *plausible* molecular graphs?

	[8]	[F]	[= 0]	[#N]	[0]	[<i>N</i>]	[= N]	[<i>C</i>]	[= <i>C</i>]	[#C]	[Ring]	[Branch1]	[Branch2]	[Branch3]
Z_0	Z_0	$\mathbf{F}Z_1$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{O}Z_2$	NZ_3	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	ign Z ₀	ign Z ₀	ign Z_0	ign Z_0
Z_1	Е	F	Ο	Ν	$\mathbf{O}Z_1$	NZ_2	NZ_2	$\mathbf{C}Z_3$	$\mathbf{C}Z_3$	$\mathbf{C}Z_3$	R(N)	ign Z ₁	ign Z_1	ign Z ₁
<i>Z</i> ₂	Е	F	=0	=N	$\mathbf{O}Z_1$	NZ_2	$=NZ_1$	$\mathbf{C}Z_3$	=C Z ₂	=C <i>Z</i> ₂	$R(N)Z_1$	$B(N, Z_5)Z_1$	$B(N, Z_5)Z_1$	$B(N, Z_5)Z_1$
Z_3	Е	F	=0	#N	$\mathbf{O}Z_1$	NZ_2	$=NZ_1$	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	#C Z ₁	$R(N)Z_2$	$B(N, Z_5)Z_2$	$B(N, Z_6)Z_1$	$B(N, Z_5)Z_2$
Z_4	Е	F	=0	#N	$\mathbf{O}Z_1$	NZ_2	$=NZ_1$	$\mathbf{C}Z_3$	=C Z ₂	#C Z ₁	$R(N)Z_3$	$B(N, Z_5)Z_3$	$B(N,Z_7)Z_1$	$B(N, Z_6)Z_3$
Z_5	С	F	Ο	Ν	$\mathbf{O}Z_1$	NZ_2	NZ_2	$\mathbf{C}Z_3$	$\mathbf{C}Z_3$	$\mathbf{C}Z_3$	Z_5	Z_5	Z_5	Z_5
Z_6	С	F	=0	=N	$\mathbf{O}Z_1$	NZ_2	$=NZ_1$	$\mathbf{C}Z_3$	=C Z ₂	=C <i>Z</i> ₂	Z_6	Z ₆	Z_6	Z_6
Z_7	С	F	=0	#N	$\mathbf{O}Z_1$	NZ_2	$=NZ_1$	$\mathbf{C}Z_3$	=C Z ₂	#C Z ₁	Z_7	Z ₇	Z_7	Z_7
Ν	1	2	3	4	5	6	7	8	9	10	11	12	13	14

// This set of decoding rules covers the generation of all *non-ionic* molecules of the *QM9* dataset



Hands-on exercise with SELFIES



	[8]	[F]	[= 0]	[#N]	[0]	[<i>N</i>]	[= N]	[<i>C</i>]	[= <i>C</i>]	[#C]	[Ring]	[Branch1]	[Branch2]	[Branch3]
Z_0	Z_0	$\mathbf{F}Z_1$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{O}Z_2$	NZ_3	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	ign Z ₀	ign Z ₀	ign Z_0	ign Z ₀
Z_1	Е	F	Ο	Ν	$\mathbf{O}Z_1$	NZ_2	NZ_2	$\mathbf{C}Z_3$	$\mathbf{C}Z_3$	$\mathbf{C}Z_3$	R(N)	ign Z ₁	ign Z ₁	ign Z_1
Z_2	Е	F	=0	=N	$\mathbf{O}Z_1$	NZ_2	$=NZ_1$	$\mathbf{C}Z_3$	=C Z ₂	=C <i>Z</i> ₂	$R(N)Z_1$	$B(N, Z_5)Z_1$	$B(N, Z_5)Z_1$	$B(N, Z_5)Z_1$
Z_3	Е	F	=0	#N	$\mathbf{O}Z_1$	NZ_2	$=NZ_1$	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	#C Z ₁	$R(N)Z_2$	$B(N, Z_5)Z_2$	$B(N,Z_6)Z_1$	$B(N, Z_5)Z_2$
Z_4	Е	F	=0	#N	$\mathbf{O}Z_1$	NZ_2	$=NZ_1$	$\mathbf{C}Z_3$	=C Z ₂	#C Z ₁	$R(N)Z_3$	$B(N, Z_5)Z_3$	$B(N,Z_7)Z_1$	$B(N,Z_6)Z_3$
Z_5	С	F	Ο	Ν	$\mathbf{O}Z_1$	NZ_2	NZ_2	$\mathbf{C}Z_3$	$\mathbf{C}Z_3$	C <i>Z</i> ₃	Z_5	Z_5	Z_5	Z_5
Z_6	С	F	=0	=N	$\mathbf{O}Z_1$	NZ_2	$=NZ_1$	$\mathbf{C}Z_3$	=C Z ₂	=C <i>Z</i> ₂	Z_6	Z_6	Z_6	Z_6
Z_7	С	F	=0	#N	$\mathbf{O}Z_1$	NZ_2	$=NZ_1$	$\mathbf{C}Z_3$	=C Z ₂	#C Z ₁	Z_7	Z_7	Z_7	Z_7
Ν	1	2	3	4	5	6	7	8	9	10	11	12	13	14

Task: Decode the given SELFIES string using (a subset of) the provided set of decoding rules!

Exercise - Decoding a SELFIES string (with selfies v0.1.1)

[O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]

Curre Z_{c}	ent sta	nte:	Cu	irrent r [O]	ule syr	nbol:	Decode OZ	d string: 2		
Deco	dingr	ule:					Decode	d SMILES		
0]] →	$\mathbf{O}Z_2$					0			
	[8]	[= 0]	[0]	[<i>N</i>]	[C]	[= <i>C</i>]	[Ring]	[Branch2]	Decoded molecule	
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign Z_0		
Z_1	ε	Ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	$R(N) Z_1$	ign Z ₁		
Z_2	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$		
Z_3	Е	=0	$\mathbf{O}Z_1$	NZ_2	C <i>Z</i> ₃	=C Z ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$		H ₂ O
Z_5	Е	Ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	Z_5	Z_5		2
Z_6	Е	=0	O Z ₁	NZ_2	C <i>Z</i> ₃	=C Z ₂	Z_5	Z_6		
Ν	1	3	5	6	8	9	11	13		

Exercise - Decoding a SELFIES string (with selfies v0.1.1)

[O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]



Exercise - Decoding a SELFIES string (with selfies v0.1.1)

[O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]

Curre Z ₂	ent sta 2	ate:	Cu	irrent i [Brai	rule syr 1 ch2]	nbol:	Decoded string: $O=C B(N, Z_5)Z_1$				
Deco [B	dingr ranc	ule: :h2] -	→ B(N	7,Z ₅)2	Z ₁		Decode O=	d SMILES C			
	[8]	[= 0]	[0]	[<i>N</i>]	[<i>C</i>]	[= <i>C</i>]	[Ring]	[Branch2]	Decoded		
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign Z ₀			
Z_1	З	0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	$R(N) Z_1$	ign Z ₁			
Z_2	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$			
Z_3	З	=0	O Z ₁	NZ_2	C Z ₃	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$	0		
Z_5	З	0	O Z ₁	NZ_2	$\mathbf{C}Z_3$	C Z ₃	Z_5	Z_5			
Z_6	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	Z_5	Z ₆			
Ν	1	3	5	6	8	9	11	13			

Exercise - Decoding a SELFIES string (with selfies v0.1.1)

[O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]

Curre N	ent sta	ate:	Cu	irrentr [eps]	rule syr 	nbol:	Decode O=	d string: C $B(1, Z_5)Z_1$			
Deco [e]	dingr ps] -	ule: → 1					Decode O=	d SMILES C	SMILES ;		
	[8]	[= 0]	[0]	[<i>N</i>]	[<i>C</i>]	[= <i>C</i>]	[Ring]	[Branch2]	Decoded		
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign Z ₀			
Z_1	Е	0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	$R(N) Z_1$	ign Z ₁			
Z_2	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$			
Z_3	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$	0		
Z_5	З	0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	Z_5	Z_5			
Z_6	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	Z_5	Z ₆			
N	1	3	5	6	8	9	11	13			

Exercise - Decoding a SELFIES string (with selfies v0.1.1)

[O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]



Exercise - Decoding a SELFIES string (with selfies v0.1.1)

[O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]



Exercise - Decoding a SELFIES string (with selfies v0.1.1)

[O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]



Exercise - Decoding a SELFIES string (with selfies v0.1.1)

[O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]

Curre Z_3	ent sta 3	ite:	Cu	irrentr [=C]	ule syn	nbol:	Decoded string: $O=C B(1,CZ_3) OC=CZ_2$				
Deco	dingr C1 →	ule:	7				Decoded O=(d SMILES C(C)OC=C			
	- C		2					• •			
	[8]	[= 0]	[0]	[N]	[C]	[= <i>C</i>]	[Ring]	[Branch2]	Decoded		
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign Z_0			
Z_1	Е	Ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	$R(N) Z_1$	ign Z ₁			
Z_2	З	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$	0		
Z_3	З	=0	O Z ₁	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$			
Z_5	ε	Ο	O Z ₁	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	Z_5	Z_5			
Z_6	З	=0	O Z ₁	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	Z_5	Z_6	Ö		
Ν	1	3	5	6	8	9	11	13			

Exercise - Decoding a SELFIES string (with selfies v0.1.1)

[O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]

Curre Z_2	ent sta	ite:	Cu	irrentr [C]	ule syr	nbol:	Decoded string: $O=C B(1,CZ_3) OC=CCZ_3$				
Deco [C	dingr] →	ule: $\mathbf{C}Z_3$					Decoded O=0	SMILES C(C)OC=CC			
	[8]	[=0]	[0]	[N]	[<i>C</i>]	[= C]	[Ring]	[Branch2]	Decoded		
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign Z ₀	molecule		
Z_1	Е	ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	$R(N) Z_1$	ign Z ₁			
Z_2	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$			
Z_3	Е	=0	$\mathbf{O}Z_1$	NZ_2	C Z ₃	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$			
Z_5	Е	ο	O Z ₁	NZ_2	$\mathbf{C}Z_3$	C Z ₃	Z_5	Z_5			
Z_6	З	=0	O Z ₁	NZ_2	C Z ₃	=C Z ₂	Z_5	Z_6			
Ν	1	3	5	6	8	9	11	13			

Exercise - Decoding a SELFIES string (with selfies v0.1.1)

[O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]

Curre Z_3	ent sta	ite:	Cu	rrentr [=C]	ule sym	nbol:	Decoded string: $O=C B(1,CZ_3) OC=CC=CZ_2$					
Deco [=(dingr C] →	ule: ▶ =C2	Z ₂				Decoded O=0	SMILES C(C)OC=CC=C				
	[8]	[= 0]	[0]	[<i>N</i>]	[<i>C</i>]	[= <i>C</i>]	[Ring]	[Branch2]	Decoded			
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign Z ₀				
Z_1	З	ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	$R(N) Z_1$	ign Z ₁				
Z_2	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$				
Z_3	З	=0	$\mathbf{O}Z_1$	NZ_2	C Z ₃	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$				
Z_5	ε	Ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	Z_5	Z_5				
Z_6	Е	=0	O Z ₁	NZ_2	C <i>Z</i> ₃	=C Z ₂	Z_5	Z_6	0			
Ν	1	3	5	6	8	9	11	13				

Exercise - Decoding a SELFIES string (with selfies v0.1.1) [O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]

Curre Z ₂	ent sta 2	te:	Cu	irrent r [C]	ule syn	n bol :	Decoded string: $O=C B(1,CZ_3) OC=CC=CCZ_3$					
Deco [C	dingri] →	ule: $\mathbf{C}Z_3$					Decoded O=0	Decoded SMILES O=C(C)OC=CC=CC				
	[8]	[=0]	[0]	[<i>N</i>]	[<i>C</i>]	[= <i>C</i>]	[Ring]	[Branch2]	Decoded			
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign Z ₀	molecule			
Z_1	Е	Ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	$R(N) Z_1$	ign Z ₁				
Z_2	З	=0	$\mathbf{O}Z_1$	NZ_2	C <i>Z</i> ₃	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$				
Z_3	З	=0	O Z ₁	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$				
Z_5	3	ο	O Z ₁	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	Z_5	Z_5				
Z_6	З	=0	O Z ₁	NZ_2	C <i>Z</i> ₃	=C Z ₂	Z_5	Z ₆	Ŭ			
Ν	1	3	5	6	8	9	11	13				

Exercise - Decoding a SELFIES string (with selfies v0.1.1) [O] [=C] [Branch2] [eps] [C] [O] [C] [=C] [C] [=C] [C] [=C] [Ring] [N] [C] [Branch2] [eps] [=O] [O]

Curre Z_3	entsta 3	ite:	Cu	irrentr [=C]	ule syn	nbol:	Decoded string: O=C $B(1,CZ_3)$ OC=CC=CZ_2					
Deco [=(dingr C] →	ule: =C2	Z ₂				Decoded SMILES O=C(C)OC=CC=C					
	[8]	[=0]	[0]	[N]	[C]	[= C]	[Ring]	[Branch2]	Decoded			
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign Z ₀				
Z_1	З	ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	$R(N) Z_1$	ign Z ₁				
Z_2	З	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$				
Z_3	З	=0	O Z ₁	NZ_2	C <i>Z</i> ₃	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$				
Z_5	ε	ο	O Z ₁	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	Z_5	Z_5				
Z_6	З	=0	O Z ₁	NZ_2	C <i>Z</i> ₃	=C Z ₂	Z_5	Z_6				
Ν	1	3	5	6	8	9	11	13				

Current state:Current rule symbol:Z2[Ring]							Decoded string: $O=C B(1,CZ_3) OC=CC=CC=CR(N)Z_2$				
Deco [R	dingr ing]	ule: $\rightarrow R$	$(N)Z_2$				Decoded SMILES O=C(C)OC=CC=CC=C				
	[8]	[=0]	[0]	[<i>N</i>]	[C]	[= C]	[Ring]	[Branch2]	Decoded		
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign ${Z}_0$			
Z_1	З	ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	$R(N) Z_1$	ign Z ₁			
Z_2	З	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$			
Z_3	З	=0	O Z ₁	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$			
Z_5	Е	ο	O Z ₁	NZ_2	$\mathbf{C}Z_3$	C Z ₃	Z_5	Z_5			
Z_6	З	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	Z_5	Z_6			
Ν	1	3	5	6	8	9	11	13			

Curre N	ent state:	Cı	urrentr [N]	ule syr	nbol:	Decoded string: O=C $B(1,CZ_3)$ OC=CC=CC=C $R(6)Z_2$				
Deco [N	dingrule:] → 6					Decoded SMILES O=C(C)OC1=CC=CC=C1				
	[8] [=	0] [0]	[<i>N</i>]	[C]	[= <i>C</i>]	[Ring]	[Branch2]	Decoded		
Z_0	<i>Z</i> ₀ O	Z_2 O Z_2	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign Z ₀	molecule		
Z_1	εΟ	$\mathbf{O}Z_1$	$\mathbf{N}Z_2$	$\mathbf{C}Z_3$	C Z ₃	$R(N) Z_1$	ign Z ₁			
Z_2	ε =0	$\mathbf{O}Z_1$	$\mathbf{N}Z_2$	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$	$\frac{1}{1}$		
Z_3	ε =0	O Z ₁	$\mathbf{N}Z_2$	C Z ₃	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$			
Z_5	εΟ	O Z ₁	NZ_2	$\mathbf{C}Z_3$	C Z ₃	Z_5	Z_5			
Z_6	ε =C	O Z ₁	$\mathbf{N}Z_2$	C Z ₃	=C Z ₂	Z_5	Z_6			
N	1 3	5	6	8	9	11	13			

	Current state:Current rule symbol:Z2[C]						Decoded string: O=C $B(1,CZ_3)$ OC=CC=CC=C $R(6)$ C Z_3 Decoded SMILES				
Decodingrule: $[C] \rightarrow CZ_3$							Decoded SMILES O=C(C)OC1=CC=CC=C1C				
	[8]	[= 0]	[0]	[<i>N</i>]	[<i>C</i>]	[= C]	[Ring]	[Branch2]	Decoded molecule		
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	C Z ₄	R(N)	ign Z ₀			
Z_1	Е	ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	$R(N) Z_1$	ign Z ₁			
Z_2	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$			
Z_3	З	=0	O Z ₁	NZ_2	C Z ₃	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$			
Z_5	ε	Ο	O Z ₁	NZ_2	$\mathbf{C}Z_3$	C Z ₃	Z_5	Z_5			
Z_6	З	=0	O Z ₁	NZ_2	C Z ₃	=C Z ₂	Z_5	Z_6			
Ν	1	3	5	6	8	9	11	13	I		

Curre Z ₃	Current state:Current rule symbol:Z3[Branch2]					nbol:	Decoded string: O=C $B(1,CZ_3)$ OC=CC=CC=C $R(6)$ C $B(N,Z_6)Z_1$				
Deco [B	dingr ranc	rule: ch2] -	→ B(N	7,Z ₆)Z	1		Decoded SMILES O=C(C)OC1=CC=CC=C1C				
	[8]	[= 0]] [0]	[<i>N</i>]	[<i>C</i>]	[=C]	[Ring]	[Branch2]	Decoded		
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	C Z ₄	R(N)	ign Z ₀			
Z_1	Е	Ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	$R(N) Z_1$	ign Z ₁			
Z_2	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$	Ĭ ſ 🕥		
Z_3	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$			
Z_5	Е	Ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	Z_5	Z_5			
Z_6	З	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	Z_5	Z_6			
Ν	1	3	5	6	8	9	11	13	I		

Curre N	Current state:Current rule symbol:N[eps]						Decoded string: O=C $B(1,CZ_3)$ OC=CC=CC=C $R(6)$ C $B(1,Z_6)Z_1$				
Decodingrule: [eps] → 1							Decoded SMILES O=C(C)OC1=CC=CC=C1C				
	[8]	[=0]	[0]	[<i>N</i>]	[C]	[= C]	[Ring]	[Branch2]	Decoded		
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	$\mathbf{C}Z_4$	R(N)	ign Z ₀	molecule		
Z_1	З	ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	$R(N) Z_1$	ign Z ₁			
Z_2	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$			
Z_3	З	=0	O Z ₁	NZ_2	C Z ₃	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$			
Z_5	З	Ο	O Z ₁	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	Z_5	Z_5			
Z_6	Е	=0	$\mathbf{O}Z_1$	NZ_2	C Z ₃	=C Z ₂	Z_5	Z ₆			
Ν	1	3	5	6	8	9	11	13	1		

Curre Z _e Deco	Current state:Current rule symbol: Z_6 [=O]Decoding rule:[=O] \rightarrow O						Decoded string: $O=C B(1,CZ_3) OC=CC=CC=C R(6) C B(1,O)Z_1$ Decoded SMILES O=C(C)OC1=CC=CC=C1C(=O)				
	[8]	[= 0]	[0]	[<i>N</i>]	[<i>C</i>]	[=C]	[Ring]	[Branch2]	Decoded		
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	CZ4	$\mathbf{C}Z_4$	R(N)	ign Z ₀			
Z_1	Е	ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	$R(N) Z_1$	ign Z ₁	\prod		
Z_2	Е	=0	$\mathbf{O}Z_1$	$\mathbf{N}Z_2$	$\mathbf{C}Z_3$	=C <i>Z</i> ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$			
Z_3	Е	=0	O Z ₁	$\mathbf{N}Z_2$	C <i>Z</i> ₃	=C <i>Z</i> ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$			
Z_5	Е	0	O Z ₁	NZ_2	$\mathbf{C}Z_3$	C <i>Z</i> ₃	Z_5	Z_5			
Z_6	Е	=O	O Z ₁	NZ_2	C Z ₃	=C Z ₂	Z_5	Z_6	0		
Ν	1	3	5	6	8	9	11	13			

Current state:Current rule symbol: Z_1 [O]Decoding rule:[O] $\rightarrow OZ_1$						nbol:	Decoded string: $O=C B(1,CZ_3) OC=CC=CC=C R(6) C B(1,O) OZ_1$ Decoded SMILES O=C(C)OC1=CC=CC=C1C(=O)O				
	[8]	[= 0]	[0]	[<i>N</i>]	[<i>C</i>]	[= <i>C</i>]	[Ring]	[Branch2]	Decoded		
Z_0	Z_0	$\mathbf{O}Z_2$	$\mathbf{O}Z_2$	NZ_3	$\mathbf{C}Z_4$	CZ4	R(N)	ign Z ₀	molecule		
Z_1	ε	ο	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	$R(N) Z_1$	ign Z ₁	<u>o</u> (\ \		
Z_2	Е	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	$R(N)Z_2$	$B(N,Z_5)Z_1$			
Z_3	З	=0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	$R(N) Z_3$	$B(N,Z_6)Z_1$	0		
Z_5	З	0	$\mathbf{O}Z_1$	NZ_2	$\mathbf{C}Z_3$	C Z ₃	Z_5	Z_5			
Z_6	Е	=0	O Z ₁	NZ_2	$\mathbf{C}Z_3$	=C Z ₂	Z_5	Z_6	acetylsalicylic HO		
Ν	1	3	5	6	8	9	11	13	acid		



Comparing SELFIES to other text-based representations



Summarizing the capabilities and limitations of SELFIES



SELFIES in comparison to other text-based representations





SELFIES in comparison to other text-based representations



Modeling opportunities with SELFIES

Which tasks could benefit from SELFIES' robustness?

Compact structural domains

Latent space of SMILES-based autoencoder



Latent space of SELFIES-based autoencoder







Modeling opportunities and challenges



Beyond generative models

42

Superfast traversal, optimization, novelty, exploration and discovery (STONED) algorithm for molecules





Beyond generative models

43

Superfast traversal, optimization, novelty, exploration and discovery (STONED) algorithm for molecules

Table 1 Number and percentage of unique molecules obtained within different fingerprint-based similarity thresholds (ô) of the starting structures. The molecules in each experiment were generated from 250 000 random string mutations of the starting structures. Additionally, for celecoxib, we also formed the local chemical space with a scaffold constraint

			Number of molecules	(and percentage)		
	Starting structure (method)	Fingerprint	$\delta > 0.75$	$\delta > 0.60$	$\delta > 0.40$	
	Aripirazole (SELFIES, random)	ECFP4	513 (0.25%)	4206 (2.15%)	34 416 (17.66%)	
	Albuterol (SELFIES, random)	FCFP4	587 (0.32%)	4156 (2.33%)	16 977 (9.35%)	
Γ	Celecoxib (SELFIES, random)	AP ECEP4	<u>478 (0.22%)</u> 198 (0.10%)	4079 (1.90%) 1925 (1.00%)	45 594 (21.66%)	
	Celecoxib (SELFIES, terminal 10%)	ECFP4	864 (2.02%)	9407 (21.99%)	34 187 (79.91%)	JELLIEJ
	Celecoxib (SELFIES, central 10%)	ECFP4	111 (0.08%)	1767 (1.32%)	15 348 (11.45%)	
Г	Celecoxib (SMILES, random)	ECFP4 ECFP4	368 (0.53%) 122 (18.43%)	7345 (10.53%) 515 (77.49%)	34 702 (49.74%) 662 (100.00%)	
L	Celecoxib (SMILES, terminal 10%)	ECFP4	90 (20.79%)	368 (84.99%)	433 (100.00%)	
	Celecoxib (SMILES, central 10%)	ECFP4	114 (22.18%)	419 (81.52%)	514 (100.00%)	SIVILES
Г	Celecoxib (SMILES, Initial 10%) Celecoxib (DeepSMILES, random)	ECFP4 ECFP4	132 (4.43%)	953 (31.99%)	2793 (93.76%)	
	Celecoxib (DeepSMILES, terminal 10%)	ECFP4	106 (9.73%)	513 (47.11%)	1083 (99.45%)	
	Celecoxib (DeepSMILES, central 10%) Celecoxib (DeepSMILES, initial 10%)	ECFP4 ECFP4	53 (6.54%) 105 (9.28%)	162 (19.98%) 609 (53.80%)	658 (81.13%) 1106 (97 70%)	Deebowiero
	Celecoxib (SELFIES, scaffold constraint)	ECFP4	354 (0.44%)	6311 (7.79%)	53 479 (66.07%)	
	Celecoxib (CReM, ChEMBL: SCScore ≤ 2.5)	ECFP4	239 (0.58%)	5547 (13.47%)	14 887 (36.14%)	

While SELFIES generates the most valid molecules, only a tiny fraction of them are structurally similar

Beyond generative models

44

Superfast traversal, optimization, novelty, exploration and discovery (STONED) algorithm for molecules



Fig. 2 Systematic local chemical space exploration of celecoxib using mutations of different SELFIES representations. The similarity is calculated using the Tanimoto distance of the ECFP4 fingerprint between celecoxib and the generated structures.



Concluding remarks



Self-referencing embedded strings as robust representation of molecules





Major contributors: Alston Lo Seyone Chithrananda

Chemistry advisor: Robert Pollice



If you want to try SELFIES, go to *https://github.com/aspuru-guzik-group/selfies* or download it with *pip install selfies We're hiring!*

Contact: florian.haese@bayer.com

a better life