

Discovering the Compositional Structure of Vector Representations with Role Learning Networks

Paul Soulos¹



Tom McCoy¹



Tal Linzen¹



Paul Smolensky^{2,1}



¹



JOHNS HOPKINS
UNIVERSITY

²



Microsoft

tl;dr

- Our technique can uncover latent compositionality in vector representations
- Interpreting compositional structure sheds light on how these models function
- We understand the inner workings well enough to write down a symbolic algorithm to produce the neural encoding
- Our approximation allows us to directly manipulate the internal representations to produce desired behavior.

What's in a compositional representation?

Consider a sequence of digits [4, 2, 7, 9]

- A set of fillers (tokens)
 - Example: {4, 2, 7, 9}

What's in a compositional representation?

Consider a sequence of digits [4, 2, 7, 9]

- A set of fillers (tokens)
 - Example: {4, 2, 7, 9}
- A set of roles (positions in the structure)
 - Example: Left-to-right {first, second, third, fourth}

What's in a compositional representation?

Consider a sequence of digits [4, 2, 7, 9]

- A set of fillers (tokens)
 - Example: {4, 2, 7, 9}
- A set of roles (positions in the structure)
 - Example: Left-to-right {first, second, third, fourth}
- A binding operation (placing a filler in a specific role filler:role)
 - Example: {4:first}

What's in a compositional representation?

Consider a sequence of digits [4, 2, 7, 9]

- A set of fillers (tokens)
 - Example: {4, 2, 7, 9}
- A set of roles (positions in the structure)
 - Example: Left-to-right {first, second, third, fourth}
- A binding operation (placing a filler in a specific role filler:role)
 - Example: {4:first}
- A composition operation (stitching all of the bound filler:roles together)
 - Example: {4:first, 2:second, 7:third, 9:fourth}

How can neural networks represent compositional structure?



How can neural networks represent compositional structure?



Tensor Product Representations (TPRs)

- A set of fillers (tokens)
- A set of roles (positions in the structure)
- A binding operation (placing a filler in a specific role filler:role)
- A composition operation (stitching all of the bound filler:roles together)

How can neural networks represent compositional structure?



Tensor Product Representations (TPRs)

- A set of fillers (tokens)

Every filler f_i is vector

- A set of roles (positions in the structure)

Every role r_i is vector

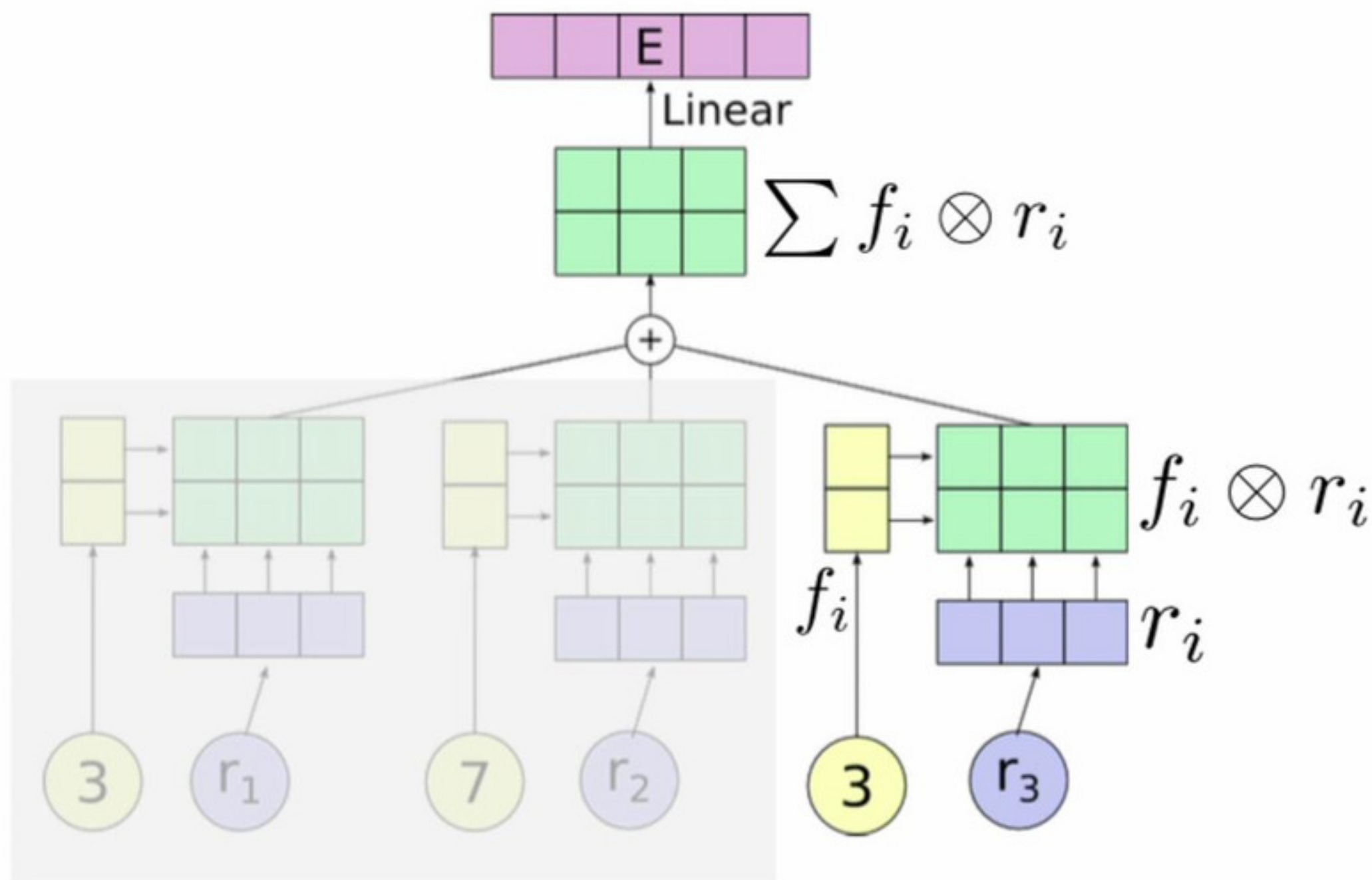
- A binding operation (placing a filler in a specific role filler:role)

Tensor product: $f_i \otimes r_i$

- A composition operation (stitching all of the bound filler:roles together)

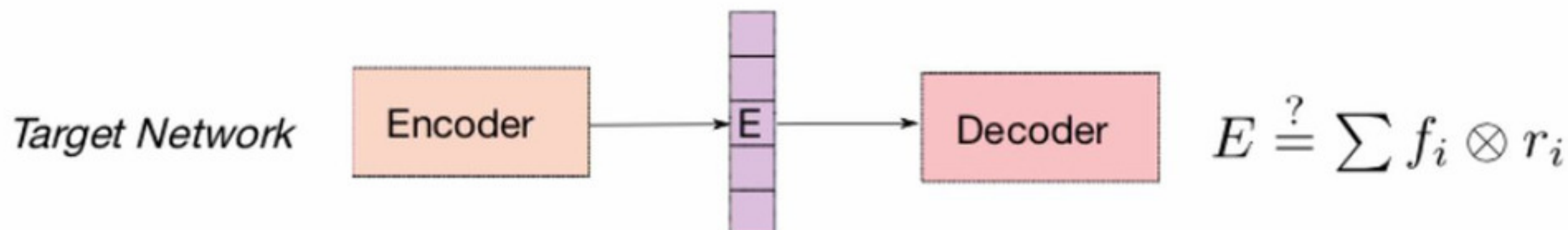
Sum: $\sum f_i \otimes r_i$

Tensor Product Encoder



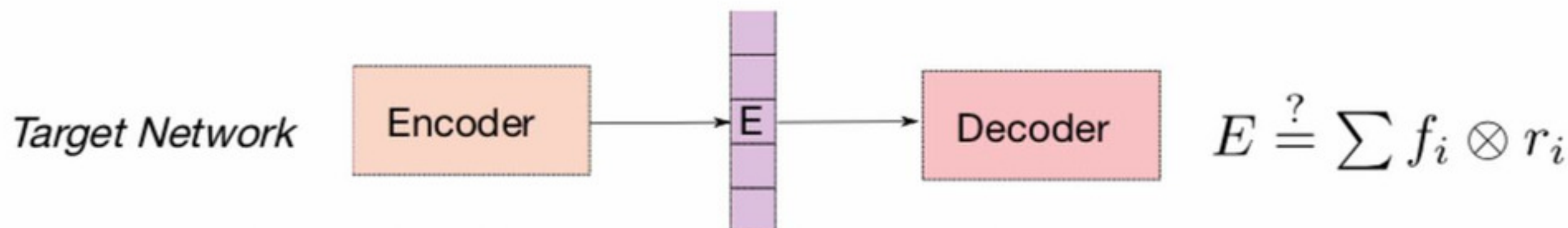
Dissecting Compositionality in Vector Representations (DISCOVER)

Goal: Discover implicit compositional structure in learned encodings E

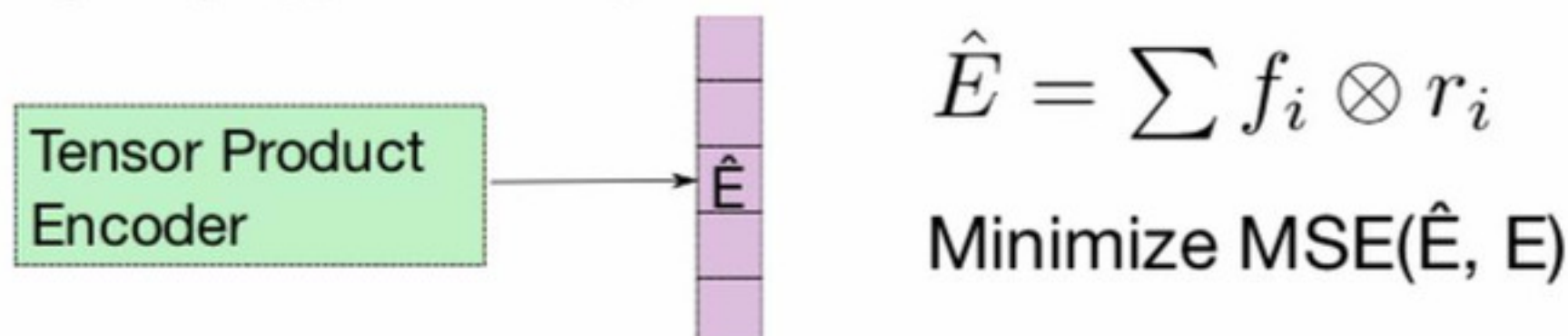


Dissecting Compositionality in Vector Representations (DISCOVER)

Goal: Discover implicit compositional structure in learned encodings E

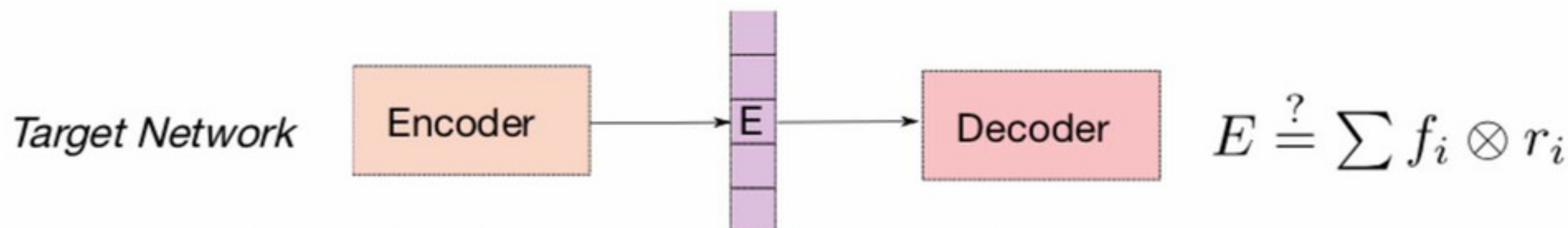


Approach: Discover implicit compositional structure in the target network's learned encoding E by approximating E with \hat{E}

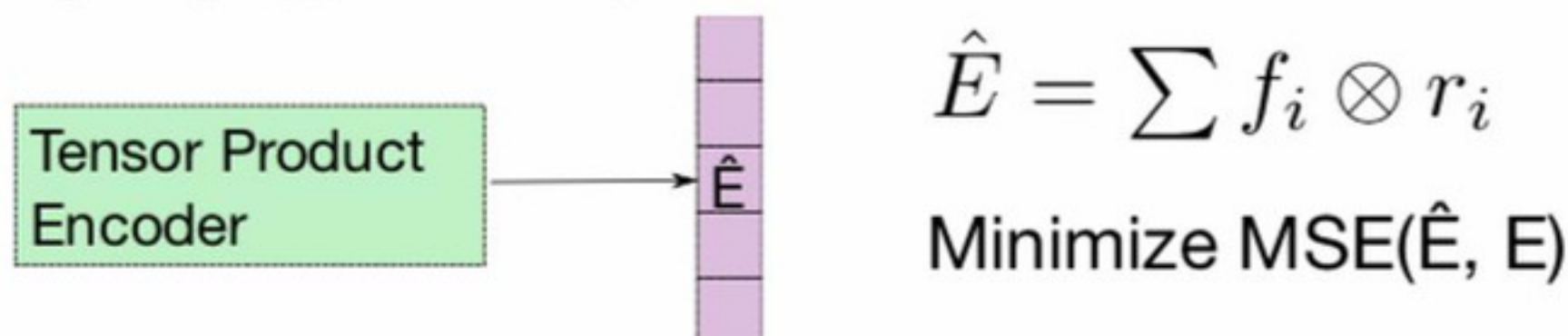


Dissecting Compositionality in Vector Representations (DISCOVER)

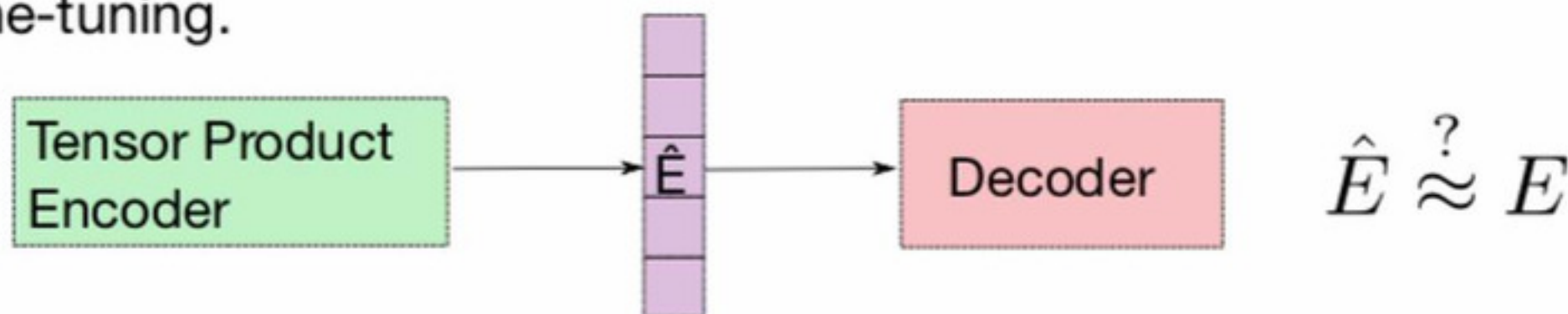
Goal: Discover implicit compositional structure in learned encodings E



Approach: Discover implicit compositional structure in the target network's learned encoding E by approximating E with \hat{E}



Evaluation: Pass the compositional encoding to the non-compositional decoder. There is no fine-tuning.

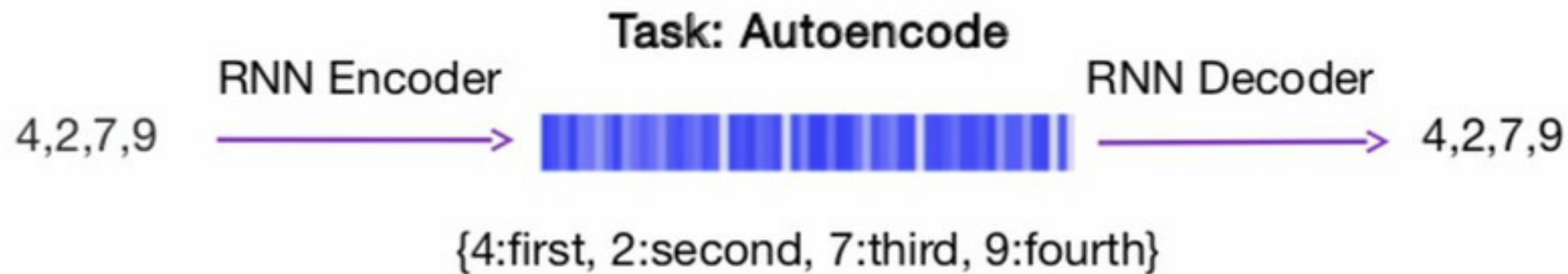


What structure is the target network learning?



Left-to-right (LTR) seems intuitive for copying. We want a FIFO queue to maintain the order.

What structure is the target network learning?

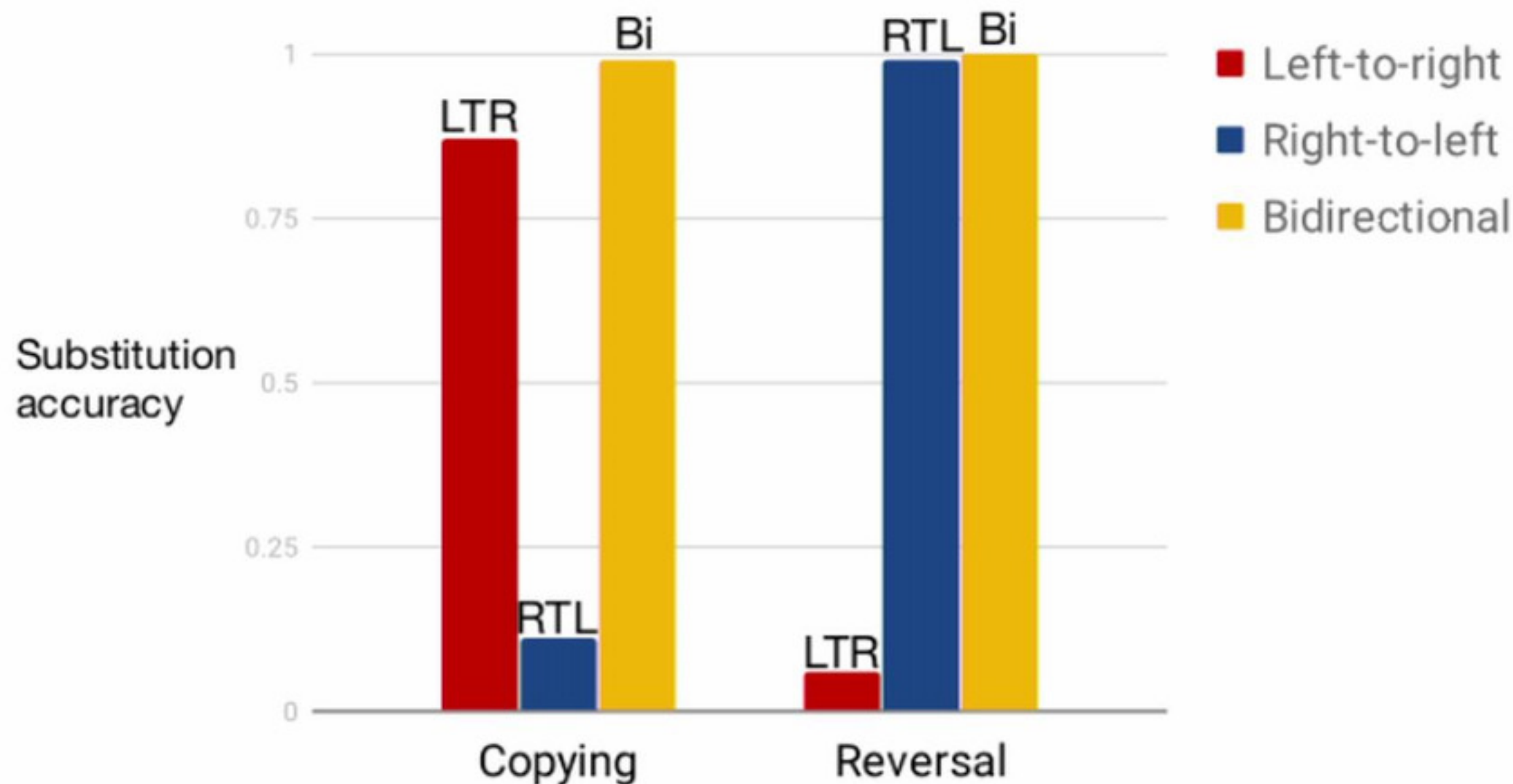


Left-to-right (LTR) seems intuitive for copying. We want a FIFO queue to maintain the order.

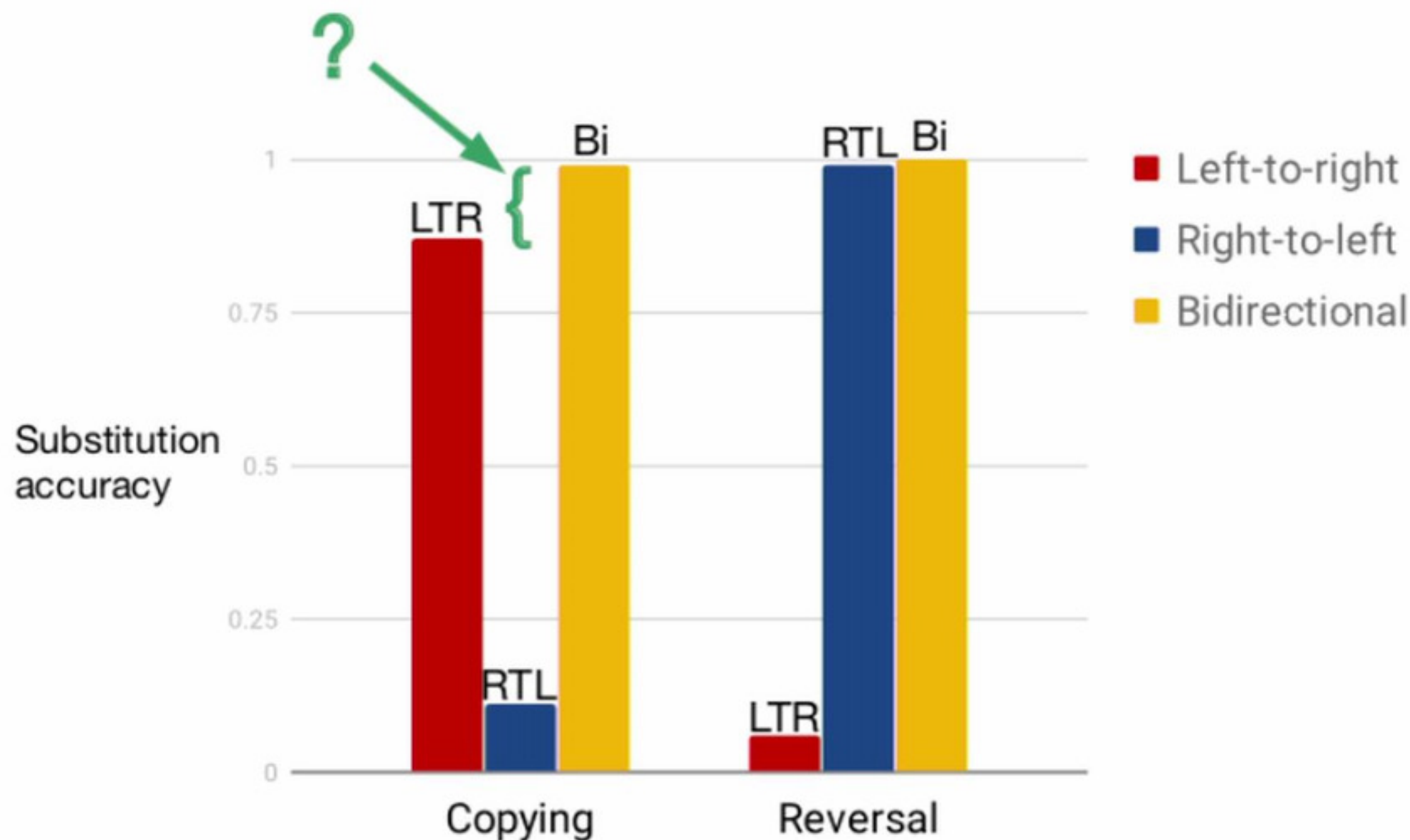


Right-to-left (RTL) seems intuitive for reversal. We want a LIFO stack to reverse the order.

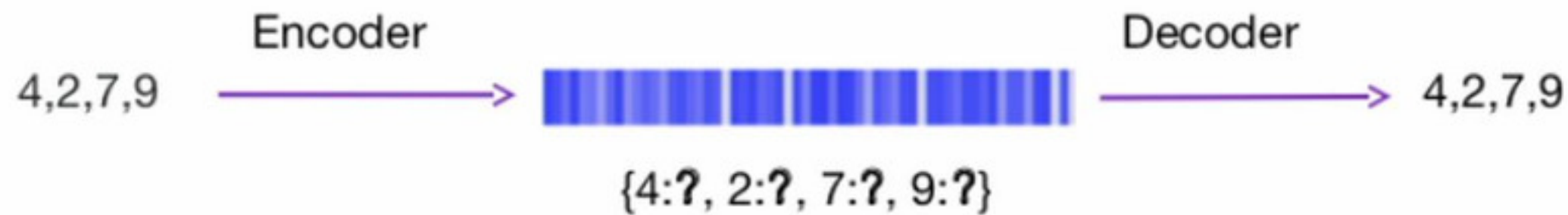
Engineered Roles



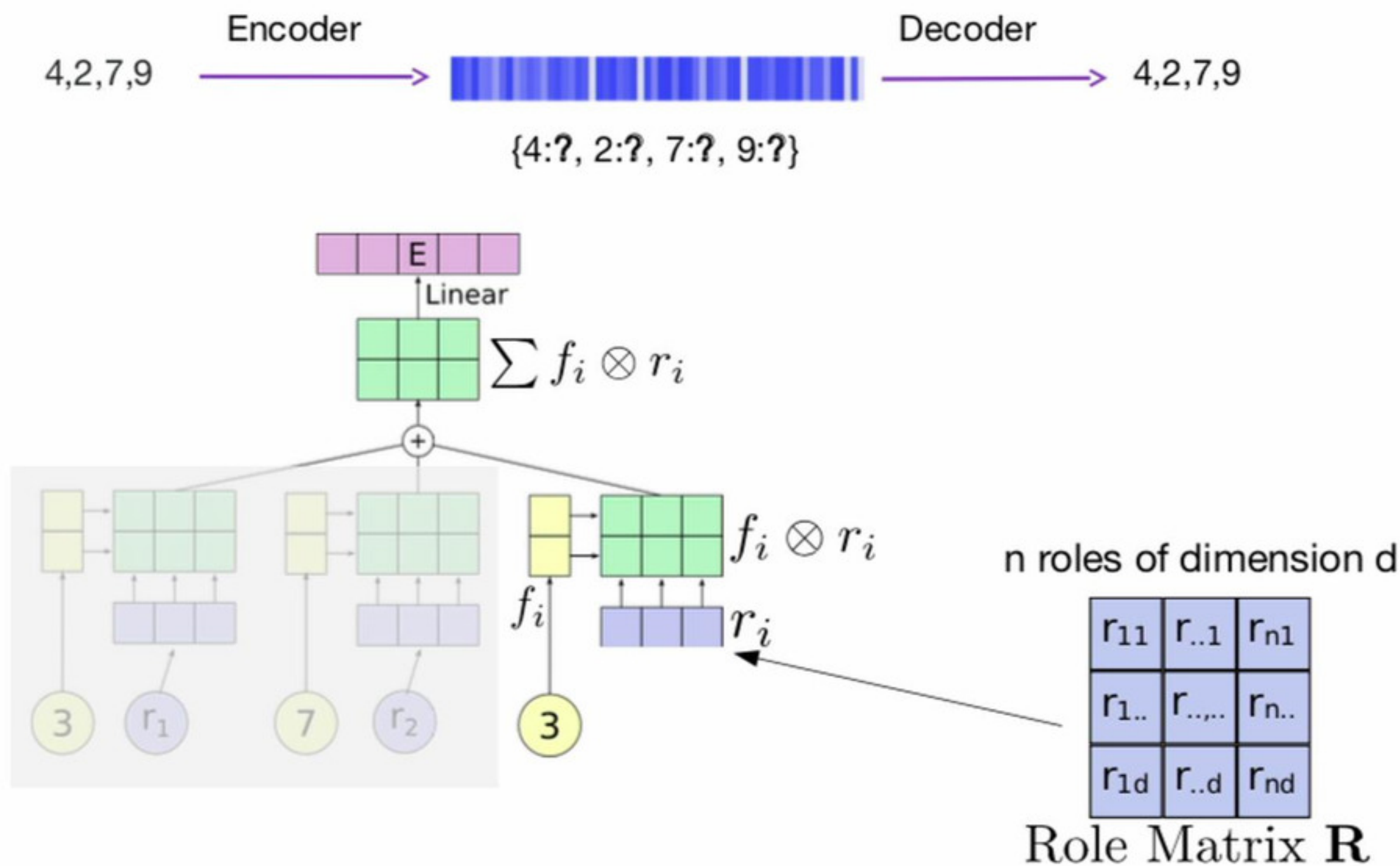
Engineered Roles



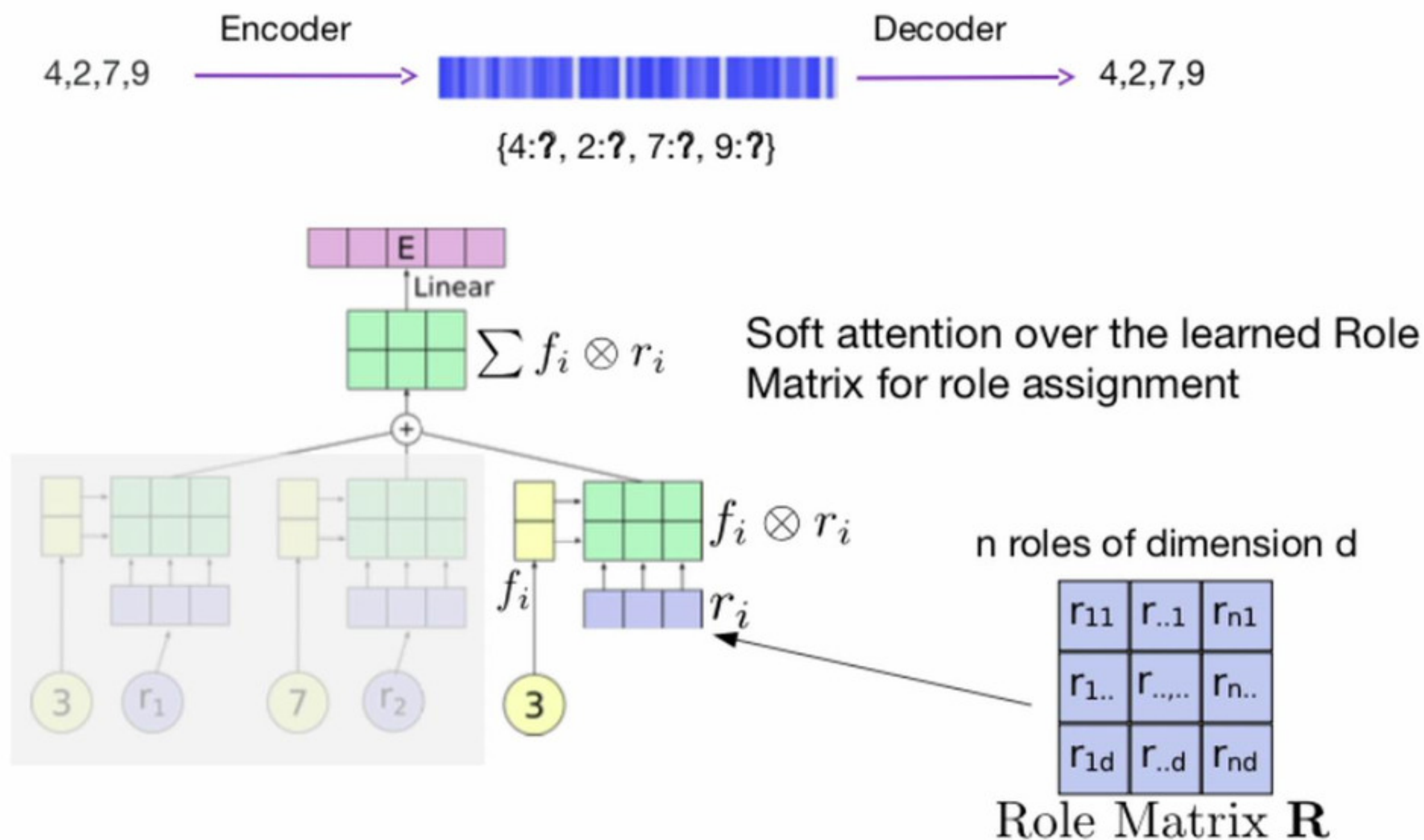
Differentiable Role Assignment



Differentiable Role Assignment



Differentiable Role Assignment



SCAN

Input	Output
jump	JUMP
jump left	LTURN JUMP
jump thrice	JUMP JUMP JUMP
jump opposite left after walk around right	RTURN WALK RTURN WALK RTURN WALK RTURN WALK LTURN LTURN JUMP

Target network is a GRU seq2seq architecture

SCAN

Input	Output
jump	JUMP
jump left	LTURN JUMP
jump thrice	JUMP JUMP JUMP
jump opposite left after walk around right	RTURN WALK RTURN WALK RTURN WALK RTURN WALK LTURN LTURN JUMP

Target network is a GRU seq2seq architecture

Target	Learned	LTR	RTL	Bi	Tree	BOW
98.5%	94.8%	6.7%	7.0%	10.7%	4.3%	4.5%

Table 1: Substitution accuracy for various encoders

SCAN Role Scheme Interpretation

- Using manual analysis of the role predictions, we created a symbolic algorithm for assigning roles to fillers
- The algorithm matches 98.7% of the role learning network's predictions on the test set.

[illegible]

SCAN

Input	Output
jump	JUMP
jump left	LTURN JUMP
jump thrice	JUMP JUMP JUMP
jump opposite left after walk around right	RTURN WALK RTURN WALK RTURN WALK RTURN WALK LTURN LTURN JUMP

Target network is a GRU seq2seq architecture

Target	Learned	LTR	RTL	Bi	Tree	BOW
98.5%	94.8%	6.7%	7.0%	10.7%	4.3%	4.5%

Table 1: Substitution accuracy for various encoders

SCAN Role Scheme Interpretation

- Using manual analysis of the role predictions, we created a symbolic algorithm for assigning roles to fillers
- The algorithm matches 98.7% of the role learning network's predictions on the test set.

```

actions = ["up", "left", "down", "right"]
actions_turn = ["up", "left", "down", "right", "turn"]
dirs = ["left", "right"]
cards = ["heuristic", "thruca"]
spins = ["opposite", "around"]

def step(sequence, correct, rplace):
    sequence = sequence.split()

    if "up" in sequence:
        part1 = sequence[sequence.index("up")+1]
        part2 = sequence[sequence.index("up")+2]
        rplace = (None for i in range(len(part2)))
        rplace[-1] = "20"

        if len(part1) == 1:
            rplace[0] = "40"
            if part1[0] == "opposite":
                if "thruca" in part1:
                    rplace[0] = "20"
                    rplace[1] = "2"
                else:
                    rplace[0] = "2"
                    rplace[1] = "4"
            elif part1[0] == "around":
                rplace[0] = "20"
                rplace[1] = "2"
            elif part1[0] in dirs:
                if len(part2) == 2:
                    part2.reverse()
                    rplace[1] = "20"
                else:
                    rplace[0] = "2"
            if len(part2) == 1:
                rplace[1] = "20"

        rplace = (None for i in range(len(part2)))
        rplace[0] = "10"
        if len(part2) == 1:
            rplace[-1] = "20"
        if len(part2) == 2:
            rplace[0] = "10"
            rplace[-1] = "2"
        if len(part2) == 3:
            rplace[0] = "10"
            rplace[1] = "20"

        return rplace + ["20"] + rplace

    elif "after" in sequence:
        part1 = sequence[sequence.index("after")+1]
        part2 = sequence[sequence.index("after")+2]
        rplace = (None for i in range(len(part2)))
        rplace[-1] = "0"

        if len(rplace) == 1:
            rplace[-1] = "20"
        if len(rplace) == 2:
            rplace[0] = "10"
        if len(rplace) == 3:
            rplace[0] = "0"
            rplace[1] = "2"

        rplace = (None for i in range(len(part2)))
        rplace[-1] = "40"

        if len(rplace) == 1:
            rplace[-1] = "4"
        if len(rplace) == 2:
            if part2[-1] == "around":
                if part2[-1] == "right":
                    rplace[0] = "4"
                elif part2[-1] == "left":
                    rplace[0] = "10"
            else:
                if part2[-1] == "thruca" and len(part2) == 2:
                    rplace[0] = "20"
                else:
                    rplace[0] = "10"
        if len(rplace) == 3:
            if part2[-1] == "thruca":
                rplace[0] = "10"
            else:
                rplace[0] = "20"

        if "10" not in rplace or (part2[-1] == "left" and part2[-1] == "around"):
            return rplace + ["10"] + rplace
        else:
            return rplace + ["40"] + rplace

    else:
        rplace = (None for i in range(len(sequence)))
        if "heuristic" in sequence or "thruca" in sequence:
            if len(sequence) == 2 and "heuristic" in sequence:
                rplace[0] = "4"
            else:
                rplace[0] = "34"

            if len(sequence) == 2 and "thruca" in sequence:
                rplace[-1] = "2"
            else:
                rplace[-1] = "40"
            if len(sequence) == 3:
                if sequence[0] in spins:
                    if sequence[0] == "opposite" and "heuristic" in sequence:
                        rplace[0] = "0"
                    elif sequence[0] == "opposite" and "thruca" in sequence:
                        rplace[0] = "34"
                    else:
                        rplace[0] = "20"
                if sequence[-1] in dirs:
                    rplace[-1] = "2"
            else:
                rplace[0] = "40"
                if len(sequence) == 1:
                    if sequence[0] == "opposite":
                        rplace[0] = "2"
                        rplace[1] = "20"
                elif sequence[0] == "around":
                    rplace[0] = "0"
                    rplace[1] = "20"
                if sequence[0] in dirs:
                    rplace[0] = "20"
                    rplace[1] = "2"

        return rplace

```

SCAN Role Scheme Interpretation

- Using manual analysis of the role predictions, we created a symbolic algorithm for assigning roles to fillers
- The algorithm matches 98.7% of the role learning network's predictions on the test set.
- Most roles are defined based on position in a subclause (e.g. *last element of the first subclause*)
- **Example roles:**
 - Role 30: Always assigned to *and*
 - Role 17: Only appears in sequences that contain the word *after*
- These two roles allow the decoder to understand the basic syntax of the command.

```
actions = ["jump", "walk", "look", "run"]
actions_turn = ["jump", "walk", "look", "run", "turn"]
dirs = ["left", "right"]
verbs = ["before", "around"]
verbs = ["before", "around"]
def assign_roles(sequence, correct_roles):
    sequence = sequence.split()
    if "and" in sequence:
        part1 = sequence[sequence.index("and")+1:]
        part2 = sequence[sequence.index("and")+1:]
        roles1 = [None for i in range(len(part1))]
        roles2 = [None for i in range(len(part2))]
        if len(part1) == 1:
            roles1[0] = "30"
        elif len(part1) == 2:
            roles1[0] = "30"
            roles1[1] = "30"
        else:
            roles1[0] = "30"
            roles1[1] = "30"
            roles1[2] = "30"
        if len(part2) == 1:
            roles2[0] = "30"
        elif len(part2) == 2:
            roles2[0] = "30"
            roles2[1] = "30"
        else:
            roles2[0] = "30"
            roles2[1] = "30"
            roles2[2] = "30"
        roles = [None for i in range(len(part1))]
        roles[0] = "30"
        if len(part1) == 1:
            roles[0] = "30"
        elif len(part1) == 2:
            roles[0] = "30"
            roles[1] = "30"
        else:
            roles[0] = "30"
            roles[1] = "30"
            roles[2] = "30"
        roles = [None for i in range(len(part2))]
        roles[0] = "30"
        if len(part2) == 1:
            roles[0] = "30"
        elif len(part2) == 2:
            roles[0] = "30"
            roles[1] = "30"
        else:
            roles[0] = "30"
            roles[1] = "30"
            roles[2] = "30"
        return roles + ["30"] + roles
    elif "after" in sequence:
        part1 = sequence[sequence.index("after")+1:]
        part2 = sequence[sequence.index("after")+1:]
        roles1 = [None for i in range(len(part1))]
        roles2 = [None for i in range(len(part2))]
        if len(part1) == 1:
            roles1[0] = "30"
        elif len(part1) == 2:
            roles1[0] = "30"
            roles1[1] = "30"
        else:
            roles1[0] = "30"
            roles1[1] = "30"
            roles1[2] = "30"
        if len(part2) == 1:
            roles2[0] = "30"
        elif len(part2) == 2:
            roles2[0] = "30"
            roles2[1] = "30"
        else:
            roles2[0] = "30"
            roles2[1] = "30"
            roles2[2] = "30"
        roles = [None for i in range(len(part1))]
        roles[0] = "30"
        if len(part1) == 1:
            roles[0] = "30"
        elif len(part1) == 2:
            roles[0] = "30"
            roles[1] = "30"
        else:
            roles[0] = "30"
            roles[1] = "30"
            roles[2] = "30"
        roles = [None for i in range(len(part2))]
        roles[0] = "30"
        if len(part2) == 1:
            roles[0] = "30"
        elif len(part2) == 2:
            roles[0] = "30"
            roles[1] = "30"
        else:
            roles[0] = "30"
            roles[1] = "30"
            roles[2] = "30"
        return roles + ["30"] + roles
    else:
        roles = [None for i in range(len(sequence))]
        if "before" in sequence or "around" in sequence:
            if len(sequence) == 2 and "before" in sequence:
                roles[0] = "30"
                roles[1] = "30"
            elif len(sequence) == 2 and "around" in sequence:
                roles[0] = "30"
                roles[1] = "30"
            else:
                roles[0] = "30"
                roles[1] = "30"
                roles[2] = "30"
        elif len(sequence) == 1:
            roles[0] = "30"
        elif len(sequence) == 2:
            roles[0] = "30"
            roles[1] = "30"
        else:
            roles[0] = "30"
            roles[1] = "30"
            roles[2] = "30"
            roles[3] = "30"
            roles[4] = "30"
            roles[5] = "30"
            roles[6] = "30"
            roles[7] = "30"
            roles[8] = "30"
            roles[9] = "30"
            roles[10] = "30"
            roles[11] = "30"
            roles[12] = "30"
            roles[13] = "30"
            roles[14] = "30"
            roles[15] = "30"
            roles[16] = "30"
            roles[17] = "30"
            roles[18] = "30"
            roles[19] = "30"
            roles[20] = "30"
            roles[21] = "30"
            roles[22] = "30"
            roles[23] = "30"
            roles[24] = "30"
            roles[25] = "30"
            roles[26] = "30"
            roles[27] = "30"
            roles[28] = "30"
            roles[29] = "30"
            roles[30] = "30"
            roles[31] = "30"
            roles[32] = "30"
            roles[33] = "30"
            roles[34] = "30"
            roles[35] = "30"
            roles[36] = "30"
            roles[37] = "30"
            roles[38] = "30"
            roles[39] = "30"
            roles[40] = "30"
            roles[41] = "30"
            roles[42] = "30"
            roles[43] = "30"
            roles[44] = "30"
            roles[45] = "30"
            roles[46] = "30"
            roles[47] = "30"
            roles[48] = "30"
            roles[49] = "30"
            roles[50] = "30"
            roles[51] = "30"
            roles[52] = "30"
            roles[53] = "30"
            roles[54] = "30"
            roles[55] = "30"
            roles[56] = "30"
            roles[57] = "30"
            roles[58] = "30"
            roles[59] = "30"
            roles[60] = "30"
            roles[61] = "30"
            roles[62] = "30"
            roles[63] = "30"
            roles[64] = "30"
            roles[65] = "30"
            roles[66] = "30"
            roles[67] = "30"
            roles[68] = "30"
            roles[69] = "30"
            roles[70] = "30"
            roles[71] = "30"
            roles[72] = "30"
            roles[73] = "30"
            roles[74] = "30"
            roles[75] = "30"
            roles[76] = "30"
            roles[77] = "30"
            roles[78] = "30"
            roles[79] = "30"
            roles[80] = "30"
            roles[81] = "30"
            roles[82] = "30"
            roles[83] = "30"
            roles[84] = "30"
            roles[85] = "30"
            roles[86] = "30"
            roles[87] = "30"
            roles[88] = "30"
            roles[89] = "30"
            roles[90] = "30"
            roles[91] = "30"
            roles[92] = "30"
            roles[93] = "30"
            roles[94] = "30"
            roles[95] = "30"
            roles[96] = "30"
            roles[97] = "30"
            roles[98] = "30"
            roles[99] = "30"
        return roles
```

Differentiable API Design

- Consider SCAN as a coding assignment between a pair of students.
 - Let's call them “Encoder” and “Decoder”

Differentiable API Design

- Consider SCAN as a coding assignment between a pair of students.
 - Let's call them “Encoder” and “Decoder”
- They split the assignment so that Encoder parses the input into a data structure, and Decoder produces the output from this data structure

? encode(List<Input Tokens>)

List<Output Tokens> decode(?)

Differentiable API Design

- Consider SCAN as a coding assignment between a pair of students.
 - Let's call them “Encoder” and “Decoder”
- They split the assignment so that Encoder parses the input into a data structure, and Decoder produces the output from this data structure

? encode(List<Input Tokens>)

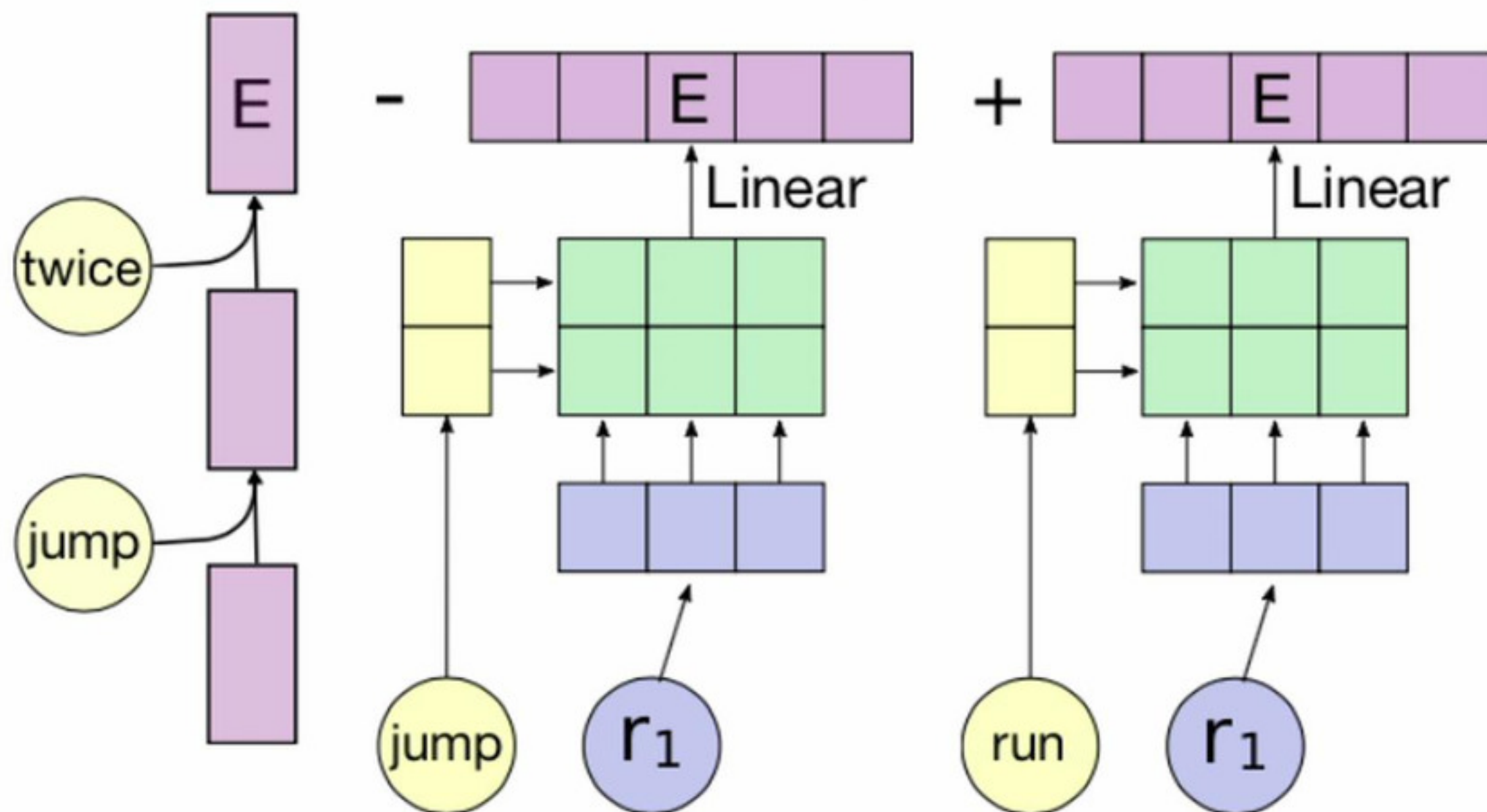
List<Output Tokens> decode(?)

<isAnd, isAfter, subclauseOneAction, subclauseOneSecondWord...> encode(List<Input Tokens>)

List<Output Tokens> decode(<isAnd, isAfter, subclauseOneAction, subclauseOneSecondWord...>)

Constituent Surgery

$$\text{emb}(\text{jump twice}) - \text{TPR}(\text{jump}) + \text{TPR}(\text{run}) \stackrel{?}{=} \text{emb}(\text{run twice})$$



JUMP JUMP → RUN RUN

Constituent Surgery

run:11 left:36 twice:8 after:43 jump:10 opposite:17 right:4 thrice:46 →
TR TR JUMP TR TR JUMP TR TR JUMP TL RUN TL RUN
– run:11 + look:11 →
TR TR JUMP TR TR JUMP TR TR JUMP TL LOOK TL LOOK
– jump:10 + walk:10 →
TR TR WALK TR TR WALK TR TR WALK TL LOOK TL LOOK



Constituent Surgery

run:11 left:36 twice:8 after:43 jump:10 opposite:17 right:4 thrice:46 →

TR TR JUMP TR TR JUMP TR TR JUMP TL RUN TL RUN

– run:11 + look:11 →

TR TR JUMP TR TR JUMP TR TR JUMP TL LOOK TL LOOK

– jump:10 + walk:10 →

TR TR WALK TR TR WALK TR TR WALK TL LOOK TL LOOK

– left:36 + right:36 →

TR TR WALK TR TR WALK TR TR WALK TR LOOK TR LOOK



Constituent Surgery

run:11 left:36 twice:8 after:43 jump:10 opposite:17 right:4 thrice:46 →

TR TR JUMP TR TR JUMP TR TR JUMP TL RUN TL RUN

– run:11 + look:11 →

TR TR JUMP TR TR JUMP TR TR JUMP TL LOOK TL LOOK

– jump:10 + walk:10 →

TR TR WALK TR TR WALK TR TR WALK TL LOOK TL LOOK

– left:36 + right:36 →

TR TR WALK TR TR WALK TR TR WALK TR LOOK TR LOOK

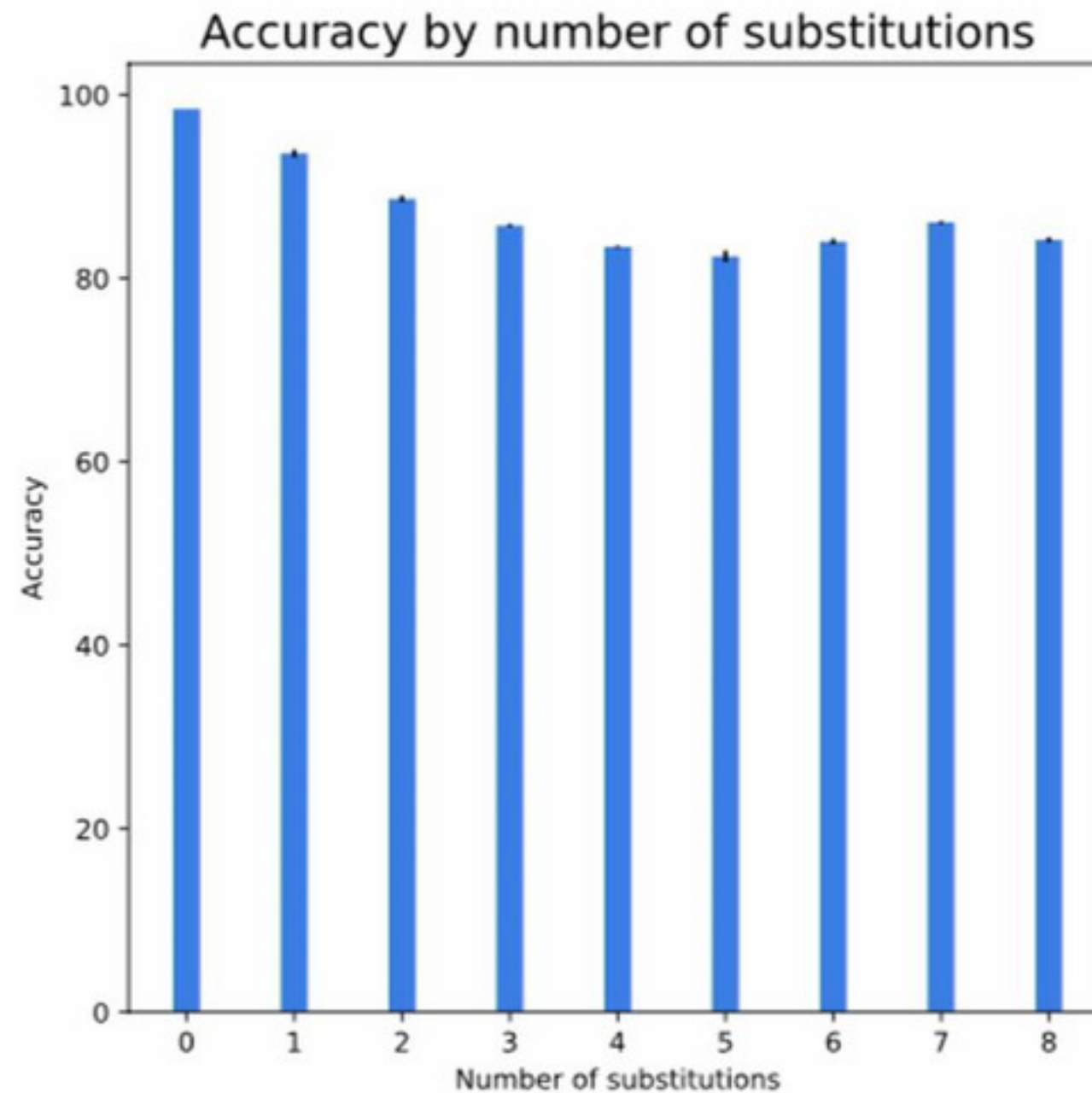
– twice:8 + thrice:8 →

TR TR WALK TR TR WALK TR TR WALK TR LOOK TR LOOK TR LOOK

Constituent Surgery

run:11 left:36 twice:8 after:43 jump:10 opposite:17 right:4 thrice:46 →
TR TR JUMP TR TR JUMP TR TR JUMP TL RUN TL RUN
– run:11 + look:11 →
TR TR JUMP TR TR JUMP TR TR JUMP TL LOOK TL LOOK
– jump:10 + walk:10 →
TR TR WALK TR TR WALK TR TR WALK TL LOOK TL LOOK
– left:36 + right:36 →
TR TR WALK TR TR WALK TR TR WALK TR LOOK TR LOOK
– twice:8 + thrice:8 →
TR TR WALK TR TR WALK TR TR WALK TR LOOK TR LOOK TR LOOK
– opposite:17 + around:17 →
TR WALK TR WALK TR WALK TR WALK TR WALK TR WALK TR WALK TR WALK TR WALK
TR WALK TR WALK TR LOOK TR LOOK TR LOOK

Constituent Surgery (Continued)



Sentence Embedding Models

	Learned	LTR	RTL	Bi	Tree	BOW
InferSent	4.05e-4	8.21e-4	9.70e-4	9.16e-4	7.78e-4	4.34e-4
Skip-thought	9.30e-5	9.91e-5	1.78e-3	3.95e-4	9.64e-5	8.87e-5
SST	5.58e-3	8.35e-3	9.29e-3	8.55e-3	5.99e-3	9.38e-3
SPINN	.139	.184	.189	.181	.178	.176

Mean-squared error for learned and engineered role schemes.

Future Directions

- Train the Tensor Product Encoder end-to-end
- Tensor Product Decoder
- Does a compositional bias improve training?
 - Train faster, fewer parameters, better generalization
- Improving natural language models with a compositional bias

Thank you!

- Run the code yourself
 - <https://github.com/psoulos/role-decomposition>
- Want more details?
 - Come by the poster
 - Check out the paper:
<https://arxiv.org/abs/1910.09113>

Acknowledgements

This material is based upon work supported by the National Science Foundation Graduate Research Fellowship Program under Grant No. 1746891, and work partially supported by NSF INSPIRE grant BCS-1344269. All opinions are our own.