# DecompX:
## Explaining Transformers Decisions by Propagating Token Decomposition
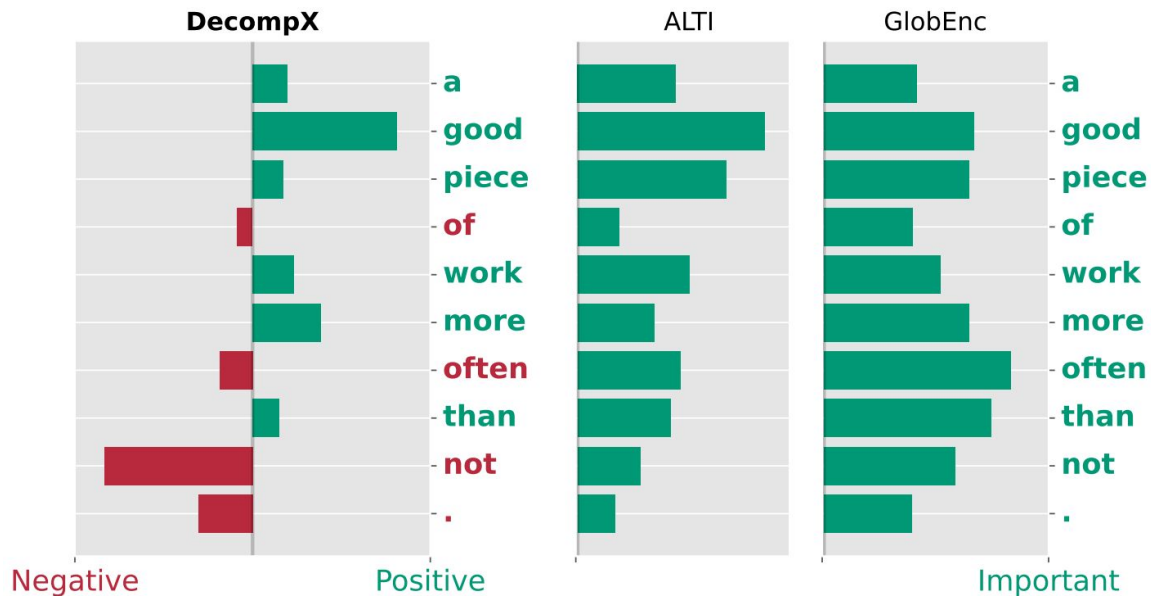
Ali Modarressi*, Mohsen Fayyaz*, Ehsan Aghazadeh,
Yadollah Yaghoobzadeh, Mohammad Taher Pilehvar
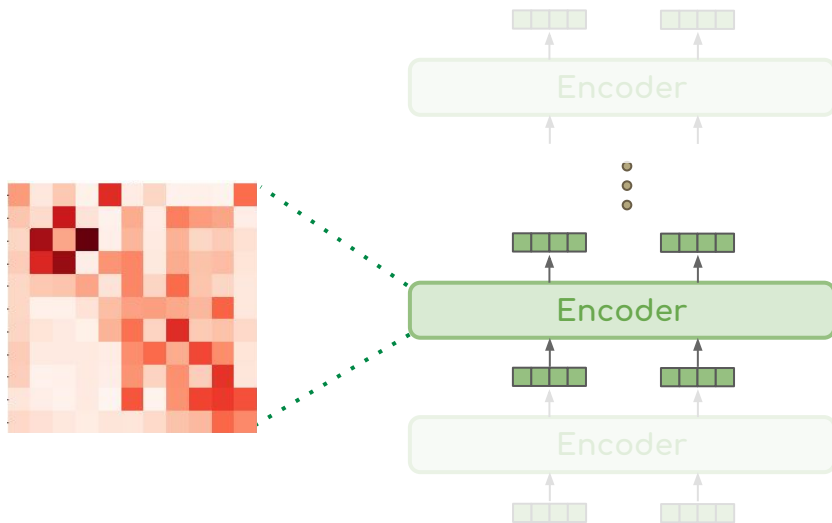
# Introduction

## What is Explanation?

[1] Ali Modarressi, Mohsen Fayyaz, Yadollah Yaghoobzadeh, and Mohammad Taher Pilehvar. 2022. GlobEnc: Quantifying global token attribution by incorporating the whole encoder layer in transformers.
[2] Javier Ferrando, Gerard I. Gállego, and Marta R. Costajussà. 2022. Measuring the mixing of contextual information in the transformer.
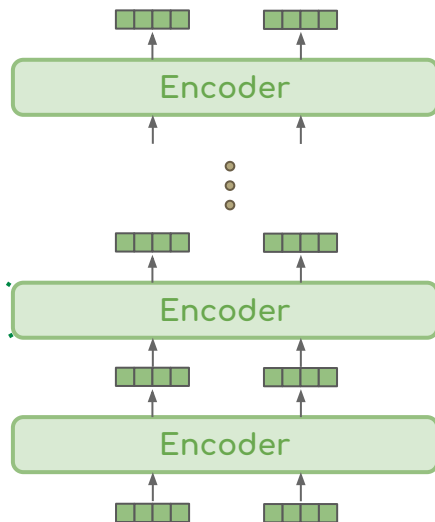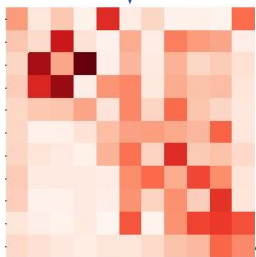
**Local Attention Map → Scalar Aggregation (e.g. Rollout, Flow)**

# Existing Methods

## Local Attention Map → Scalar Aggregation (e.g. Rollout, Flow)

- **Raw-attention[1]**
- **ALTI[2]**
- **Globenc[3]**
- **Value-Zeroing[4]**

[1] Samira Abnar and Willem Zuidema. 2020. Quantifying attention flow in transformers. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 4190–4197, Online. Association for Computational Linguistics.
[2] Javier Ferrando, Gerard I. Gállego, and Marta R. Costajussà. 2022. Measuring the mixing of contextual information in the transformer.
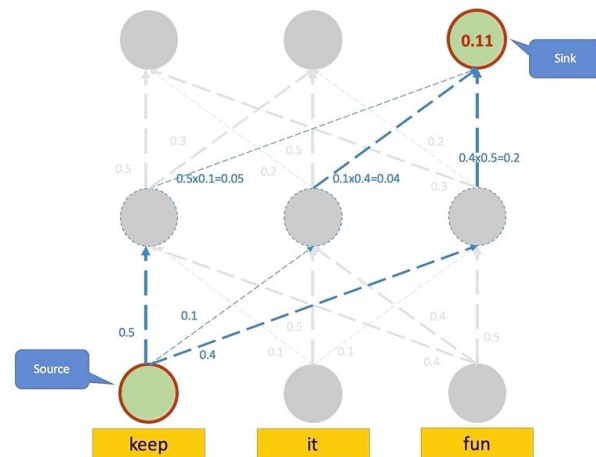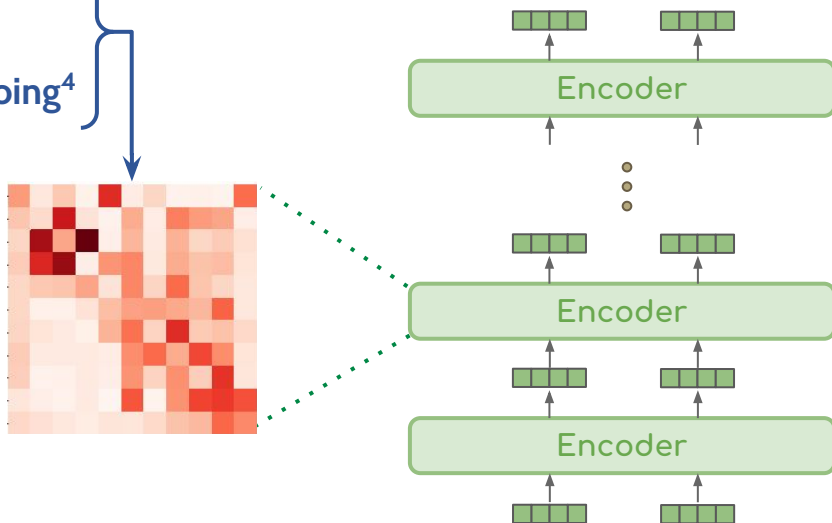[3] Ali Modarressi, Mohsen Fayyaz, Yadollah Yaghoobzadeh, and Mohammad Taher Pilehvar. 2022. GlobEnc: Quantifying global token attribution by incorporating the whole encoder layer in transformers.
[4] Hosein Mohebbi, Willem Zuidema, Grzegorz Chrupała, and Afra Alishahi. 2023. Quantifying context mixing in transformers. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics

# Existing Methods

## Local Attention Map → Scalar Aggregation (e.g. Rollout, Flow)

- **Raw-attention[1]**
- **ALTI[2]**
- **Globenc[3]**
- **Value-Zeroing[4]**

[1] Samira Abnar and Willem Zuidema. 2020. Quantifying attention flow in transformers. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 4190–4197, Online. Association for Computational Linguistics.
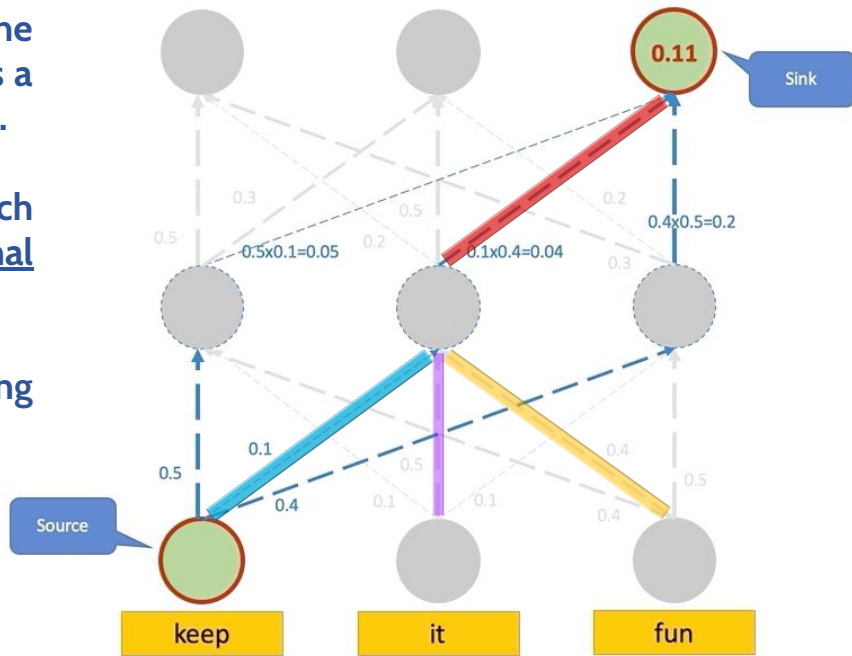[2] Javier Ferrando, Gerard I. Gállego, and Marta R. Costajussà. 2022. Measuring the mixing of contextual information in the transformer.
[3] Ali Modarressi, Mohsen Fayyaz, Yadollah Yaghoobzadeh, and Mohammad Taher Pilehvar. 2022. GlobEnc: Quantifying global token attribution by incorporating the whole encoder layer in transformers.
[4] Hosein Mohebbi, Willem Zuidema, Grzegorz Chrupała, and Afra Alishahi. 2023. Quantifying context mixing in transformers. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics
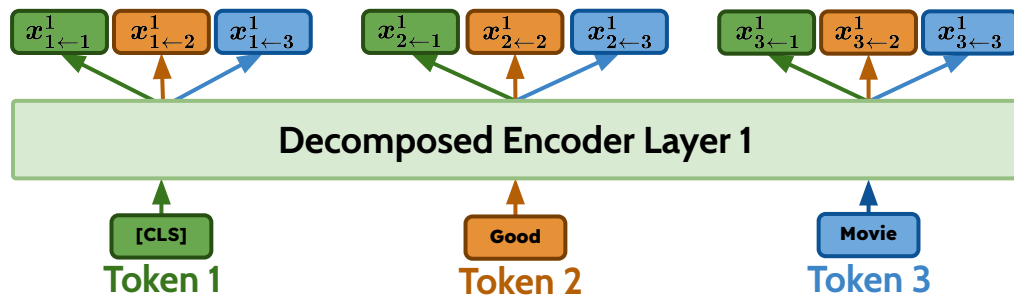
# Scalar Aggregation Issues

➜ Scalar aggregation methods (e.g. Rollout) assume that the only required information for computing the global flow is a set of <u>scalar cross-token attributions</u>.

➜ Nevertheless, this simplifying assumption ignores that each decomposed vector represents the <u>multi-dimensional impact of its inputs</u>.

➜ Therefore, losing information is inevitable when reducing these complex vectors into one cross-token weight.
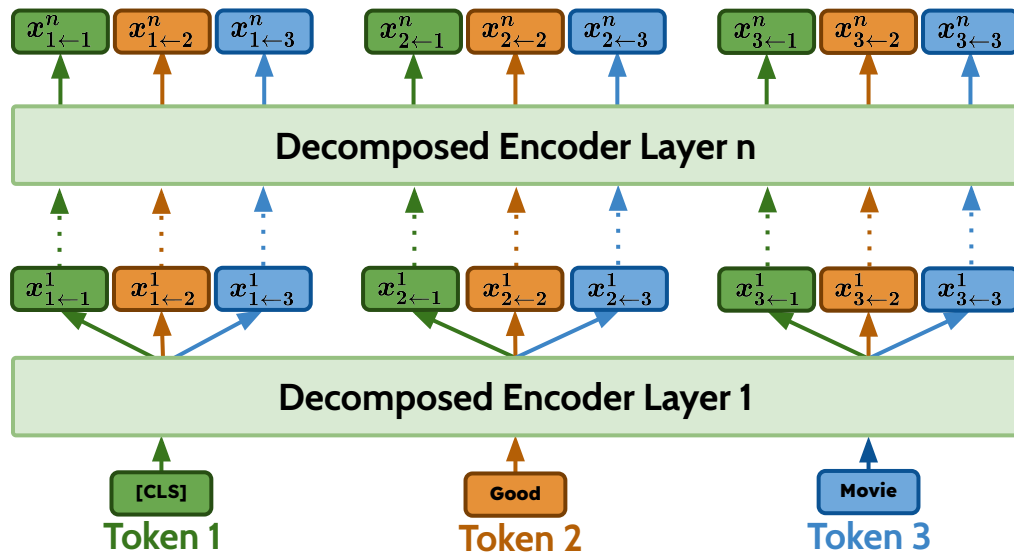
# Our Solution: DecompX

**Propagating Token Decomposition**

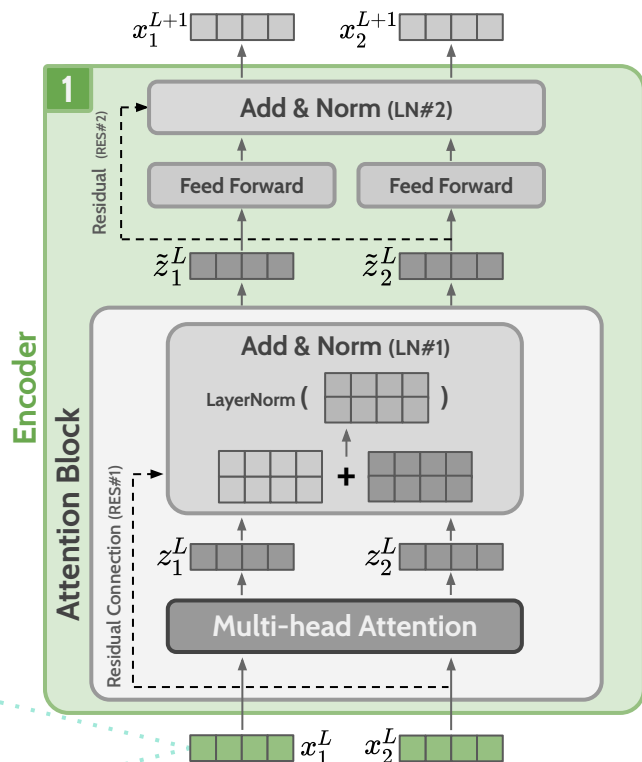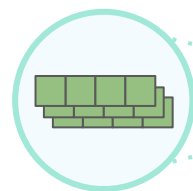# Our Solution: DecompX

## Propagating Token Decomposition

# DecompX

## Propagating Token Decomposition
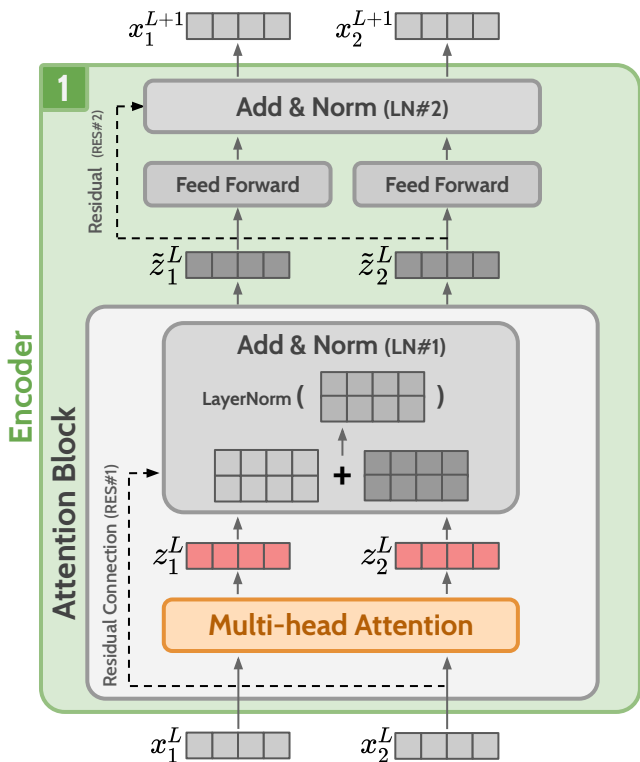
$$x_i^\ell = \sum_{k=1}^{N} x_{i \Leftarrow k}^\ell$$

# DecompX

## Propagating Token Decomposition

$$\boldsymbol{x}_{i \Leftarrow k}^{\ell}$$

$$\boldsymbol{z}_i^{\ell} = \sum_{k=1}^{N} \big( \underbrace{\sum_{h=1}^{H} \sum_{j=1}^{N} \alpha_{i,j}^h \boldsymbol{x}_{j \Leftarrow k}^{\ell} \boldsymbol{W}_{Att}^h + \omega_k \boldsymbol{b}_{Att}}_{\boldsymbol{z}_{i \Leftarrow k}^{\ell}} \big)$$

# DecompX

## Propagating Token Decomposition

$$\boldsymbol{x}_{i \Leftarrow k}^{\ell}$$

$$\boldsymbol{z}_i^{\ell} = \sum_{k=1}^{N} \left( \underbrace{\sum_{h=1}^{H} \sum_{j=1}^{N} \alpha_{i,j}^{h} \boldsymbol{x}_{j \Leftarrow k}^{\ell} \boldsymbol{W}_{\boldsymbol{Att}}^{h}}_{\boldsymbol{z}_{i \Leftarrow k}^{\ell}} + \omega_k \boldsymbol{b}_{\boldsymbol{Att}} \right)$$

$$\mathrm{LN}(\boldsymbol{z}^+{}_i^{\ell}) = \sum_{k=1}^{N} \underbrace{g_{\boldsymbol{z}^+{}_i^{\ell}}(\boldsymbol{z}^+{}_{i \Leftarrow k}^{\ell}) + \boldsymbol{\beta}}_{\tilde{\boldsymbol{z}}_{i \Leftarrow k}^{\ell}}$$

# DecompX

## Propagating Token Decomposition

$$\boldsymbol{x}_{i \Leftarrow k}^{\ell}$$

$$z_i^{\ell} = \sum_{k=1}^{N} (\underbrace{\sum_{h=1}^{H} \sum_{j=1}^{N} \alpha_{i,j}^h \boldsymbol{x}_{j \Leftarrow k}^{\ell} \boldsymbol{W}_{Att}^h}_{\boldsymbol{z}_{i \Leftarrow k}^{\ell}} + \omega_k \boldsymbol{b}_{Att})$$

$$\mathrm{LN}(\boldsymbol{z^+}_i^{\ell}) = \sum_{k=1}^{N} \underbrace{g_{\boldsymbol{z^+}_i^{\ell}}(\boldsymbol{z^+}_{i \Leftarrow k}^{\ell}) + \boldsymbol{\beta}}_{\tilde{\boldsymbol{z}}_{i \Leftarrow k}^{\ell}}$$

$$\boldsymbol{z}_{\mathrm{FFN},i}^{\ell} = f_{\mathrm{act}}^{(\boldsymbol{\zeta}_i^{\ell})} (\sum_{k=1}^{N} \boldsymbol{\zeta}_{i \Leftarrow k}^{\ell}) \boldsymbol{W}_{\mathrm{FFN}}^2 + \boldsymbol{b}_{\mathrm{FFN}}^2$$

$$= \sum_{k=1} \underbrace{\boldsymbol{\theta}^{(\boldsymbol{\zeta}_i^{\ell})} \odot \boldsymbol{\zeta}_{i \Leftarrow k}^{\ell} + \boldsymbol{b}_{\mathrm{FFN}}^2}_{\boldsymbol{z}_{\mathrm{FFN},i \Leftarrow k}^{\ell}}$$

# DecompX

## Propagating Token Decomposition

$$\boldsymbol{x}_{i \Leftarrow k}^{\ell}$$

$$z_i^{\ell} = \sum_{k=1}^{N} \left( \sum_{h=1}^{H} \sum_{j=1}^{N} \alpha_{i,j}^{h} \boldsymbol{x}_{j \Leftarrow k}^{\ell} \boldsymbol{W}_{Att}^{h} + \omega_k \boldsymbol{b}_{Att} \right)}_{z_{i \Leftarrow k}^{\ell}}$$

$$\mathrm{LN}(\boldsymbol{z^{+}}_i^{\ell}) = \sum_{k=1}^{N} \underbrace{g_{\boldsymbol{z^{+}}_i^{\ell}}(\boldsymbol{z^{+}}_{i \Leftarrow k}^{\ell}) + \boldsymbol{\beta}}_{\tilde{z}_{i \Leftarrow k}^{\ell}}$$

$$z_{\mathrm{FFN},i}^{\ell} = f_{\mathrm{act}}^{(\zeta_i^{\ell})} \left( \sum_{k=1}^{N} \zeta_{i \Leftarrow k}^{\ell} \right) \boldsymbol{W}_{\mathrm{FFN}}^{2} + \boldsymbol{b}_{\mathrm{FFN}}^{2}$$

$$= \sum_{k=1}^{N} \underbrace{\boldsymbol{\theta}^{(\zeta_i^{\ell})} \odot \zeta_{i \Leftarrow k}^{\ell} + \boldsymbol{b}_{\mathrm{FFN}}^{2}}_{z_{\mathrm{FFN},i \Leftarrow k}^{\ell}}$$

$$\boldsymbol{x}_i^{\ell+1} = \mathrm{LN}\left( \sum_{k=1}^{N} [\underbrace{\tilde{\boldsymbol{z}}_{i \Leftarrow k}^{\ell} + \boldsymbol{z}_{\mathrm{FFN},i \Leftarrow k}^{\ell}}_{\boldsymbol{z}_{\mathrm{FFN}+,i \Leftarrow k}^{\ell}}] \right)$$

$$= \sum_{k=1}^{N} \underbrace{g_{\boldsymbol{z}_{\mathrm{FFN}+,i}^{\ell}}(\boldsymbol{z}_{\mathrm{FFN}+,i \Leftarrow k}^{\ell}) + \boldsymbol{\beta}}_{\boldsymbol{x}_{i \Leftarrow k}^{\ell+1}}$$
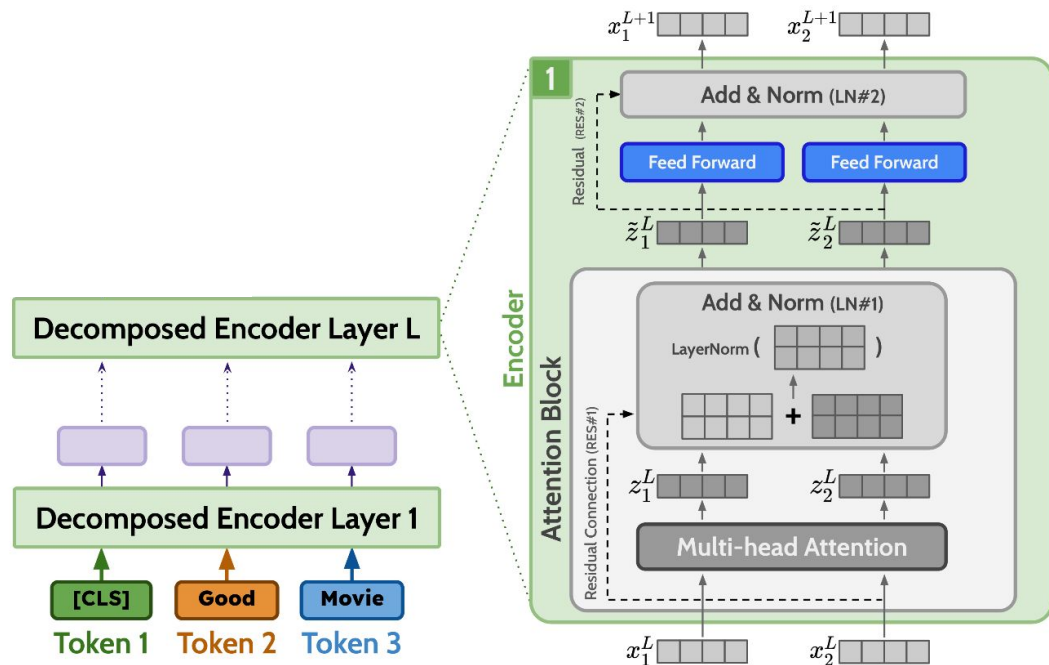
# DecompX

## Propagating Token Decomposition

# DecompX

## Propagating Token Decomposition
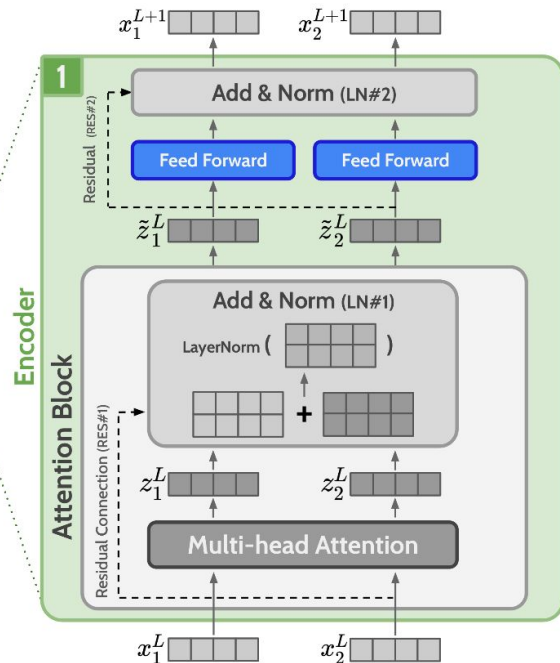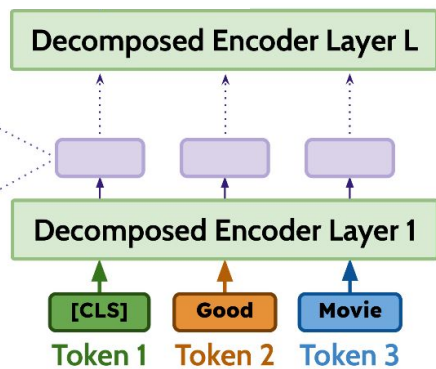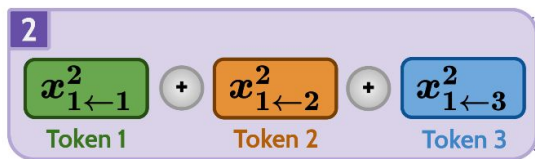
# DecompX

## Propagating Token Decomposition



$$x_{i \Leftarrow k}^{\ell} \longrightarrow x_{i \Leftarrow k}^{\ell+1}$$

**Propagate through L layers**
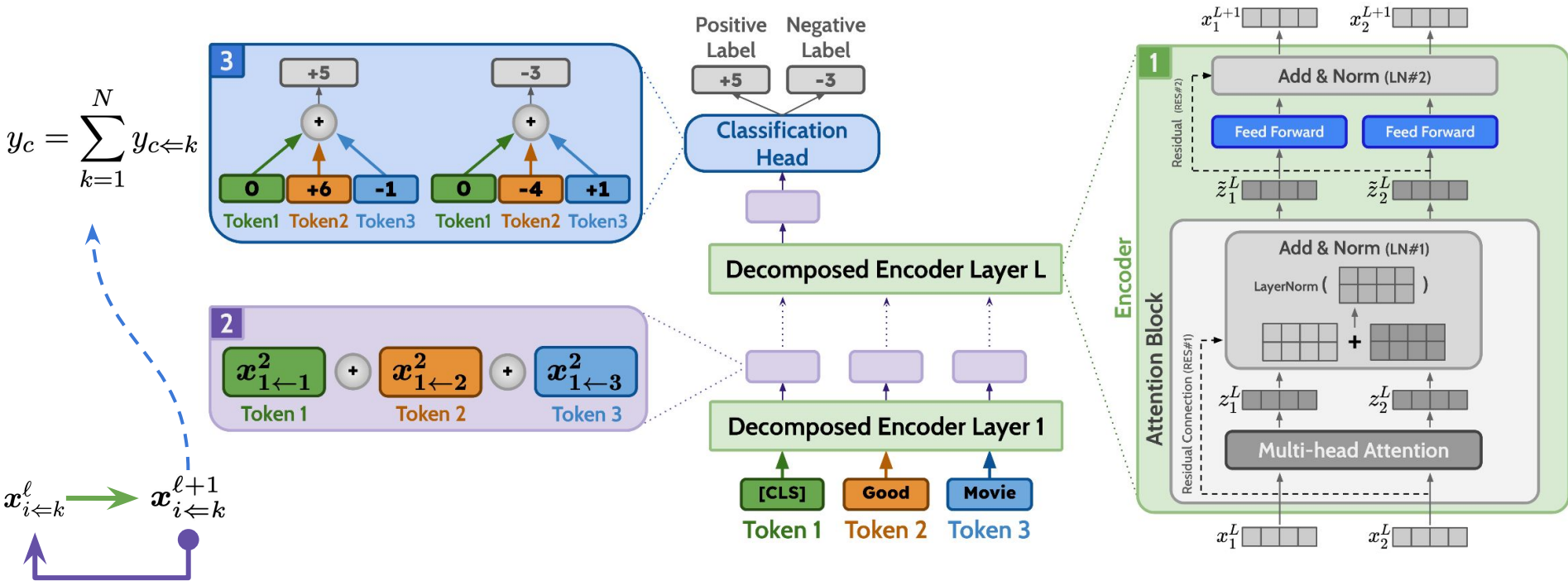
**2**

$x_{1 \leftarrow 1}^{2}$ ⊕ $x_{1 \leftarrow 2}^{2}$ ⊕ $x_{1 \leftarrow 3}^{2}$

Token 1    Token 2    Token 3

**Decomposed Encoder Layer L**

**Decomposed Encoder Layer 1**

[CLS]    Good    Movie

Token 1    Token 2    Token 3

**1**

**Encoder**

**Attention Block**

$x_1^{L+1}$    $x_2^{L+1}$

Add & Norm (LN#2)

Residual (RES#2)

Feed Forward    Feed Forward

$\tilde{z}_1^L$    $\tilde{z}_2^L$

Add & Norm (LN#1)

LayerNorm ( )

$+$

Residual Connection (RES#1)

$z_1^L$    $z_2^L$

**Multi-head Attention**

$x_1^L$    $x_2^L$

# DecompX

## Propagating Token Decomposition



$$y_c = \sum_{k=1}^{N} y_{c \Leftarrow k}$$

$$\boldsymbol{x}_{i \Leftarrow k}^{\ell} \longrightarrow \boldsymbol{x}_{i \Leftarrow k}^{\ell+1}$$

**3**

| +5 | | -3 |
|---|---|---|

| Token1 | Token2 | Token3 | Token1 | Token2 | Token3 |
|---|---|---|---|---|---|
| 0 | +6 | -1 | 0 | -4 | +1 |

**Positive Label** +5    **Negative Label** -3

**Classification Head**

**Decomposed Encoder Layer L**

**2**
$$\boldsymbol{x}_{1 \leftarrow 1}^{2} \quad \oplus \quad \boldsymbol{x}_{1 \leftarrow 2}^{2} \quad \oplus \quad \boldsymbol{x}_{1 \leftarrow 3}^{2}$$
Token 1    Token 2    Token 3

**Decomposed Encoder Layer 1**

[CLS]    Good    Movie
Token 1    Token 2    Token 3

**1**

Encoder

Attention Block

Residual (RES#2)

$x_1^{L+1}$    $x_2^{L+1}$

**Add & Norm (LN#2)**

Feed Forward    Feed Forward

$\tilde{z}_1^L$    $\tilde{z}_2^L$

**Add & Norm (LN#1)**

LayerNorm ( ⊞ )

⊞ + ⊞

Residual Connection (RES#1)

$z_1^L$    $z_2^L$

**Multi-head Attention**

$x_1^L$    $x_2^L$

# DecompX

## Overview



Our contributions:

1) We incorporated all the encoder layer components including nonlinear functions
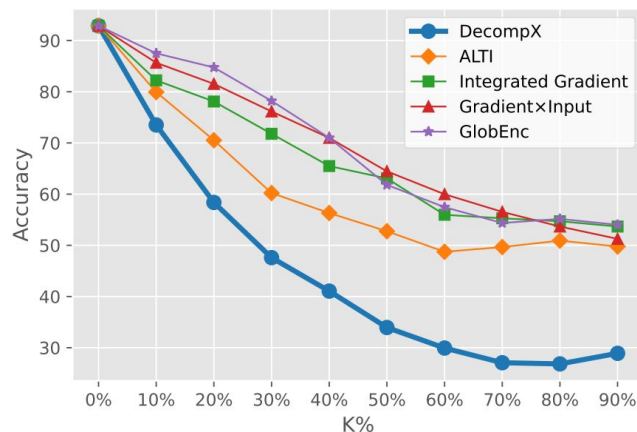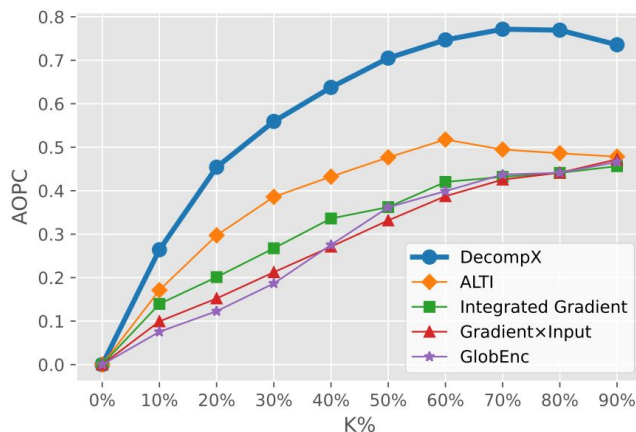
2) Propagated the decomposed vectors throughout the whole model

3) Incorporated the classification head

# Evaluation

## Results



$$\mathrm{AOPC}(K) = \frac{1}{N} \sum_{i=1}^{N} p(\hat{y} \mid x_i) - p(\hat{y} \mid \tilde{x}_i^{(K)})$$

➔ AOPC and Accuracy of different explanation methods on SST2 upon masking K% of the most important tokens.

➔ DecompX outperforms existing explanation methods, both vector- and gradient-based, by a large margin at every corruption ratio.

# Evaluation

## Aggregated Results

| | SST2 | | | MNLI | | | QNLI | | | HATEXPLAIN | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ACC↓ | AOPC↑ | PRED↑ | ACC↓ | AOPC↑ | PRED↑ | ACC↓ | AOPC↑ | PRED↑ | ACC↓ | AOPC↑ | PRED↑ |
| GlobEnc (Modarressi et al., 2022) | 67.14 | 0.307 | 72.36 | 48.07 | 0.498 | 70.43 | 64.93 | 0.342 | 84.00 | 47.65 | 0.401 | 56.50 |
| + FFN | 64.90 | 0.326 | 79.01 | 45.05 | 0.533 | 75.15 | 63.74 | 0.354 | 84.97 | 46.89 | 0.406 | 59.52 |
| ALTI (Ferrando et al., 2022) | 57.65 | 0.416 | 88.30 | 45.89 | 0.515 | 74.24 | 63.85 | 0.355 | 85.69 | 43.30 | 0.469 | 64.67 |
| Gradient×Input | 66.69 | 0.310 | 67.20 | 44.21 | 0.544 | 76.05 | 62.93 | 0.366 | 86.27 | 46.28 | 0.433 | 60.67 |
| Integrated Gradients | 64.48 | 0.340 | 64.56 | 40.80 | 0.579 | 73.94 | 61.12 | 0.381 | 86.27 | 45.19 | 0.445 | 64.46 |
| **DecompX** | **40.80** | **0.627** | **92.20** | **32.64** | **0.703** | **80.95** | **57.50** | **0.453** | **89.84** | **38.71** | **0.612** | **66.34** |

➜ **Accuracy, AOPC, and Prediction Performance of DecompX compared with the existing methods on different datasets.**

➜ **DecompX consistently outperforms other methods, which confirms that a holistic vector-based approach can present higher-quality explanations.**

# Online Demo

# THANK YOU!

amodaresi@cis.lmu.de

mohsen.fayyaz77@ut.ac.ir

Eaghazade1998@ut.ac.ir

y.yaghoobzadeh@ut.ac.ir

mp792@cam.ac.uk

Github.com/mohsenfayyaz/DecompX