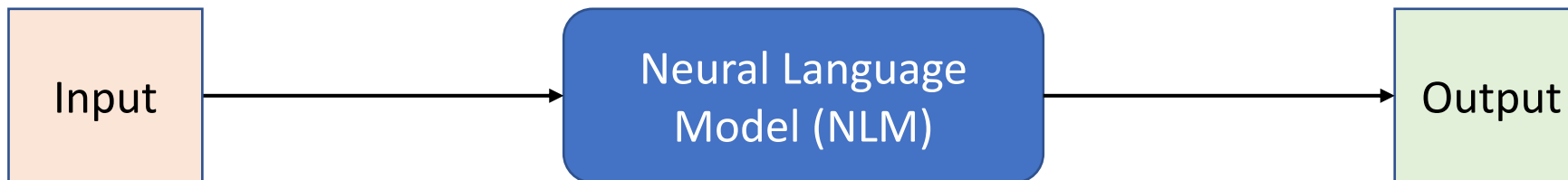# UNIREX: A Unified Learning Framework for Language Model Rationale Extraction

**Aaron Chan**, Maziar Sanjabi, Lambert Mathias, Liang Tan,

Shaoliang Nie, Xiaochang Peng, Xiang Ren, Hamed Firooz

ICML 2022

# Rationale Extraction

# Three Desiderata of Rationale Extraction



Rationale Extractor

Still , this flick is **fun** , and host to some truly **excellent** sequences .

# Three Desiderata of Rationale Extraction

**1. Faithfulness**

The explanation accurately reflects my reasoning process! 👍

Neural Language Model (NLM)

Rationale Extractor

Still , this flick is **<u>fun</u>** , and host to some truly **<u>excellent</u>** sequences .

# Three Desiderata of Rationale Extraction

**1. Faithfulness**

The explanation accurately reflects my reasoning process! 👍

Neural Language Model (NLM)

Rationale Extractor

Still , this flick is **<u>fun</u>** , and host to some truly **<u>excellent</u>** sequences .

**2. Plausibility**
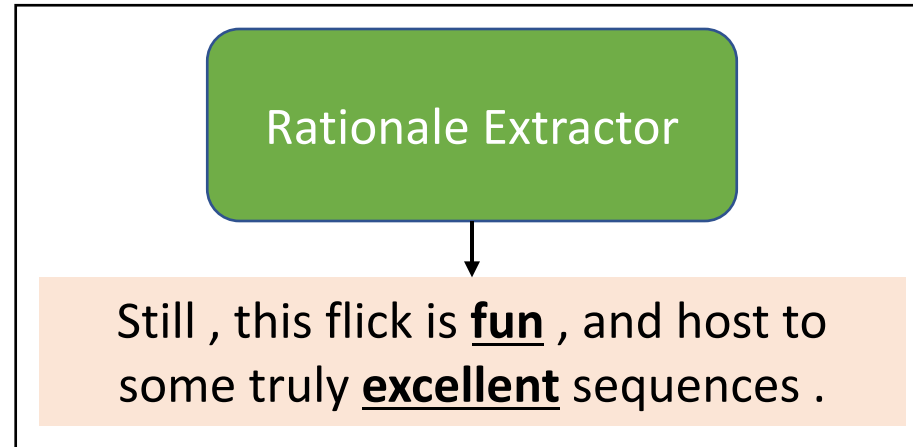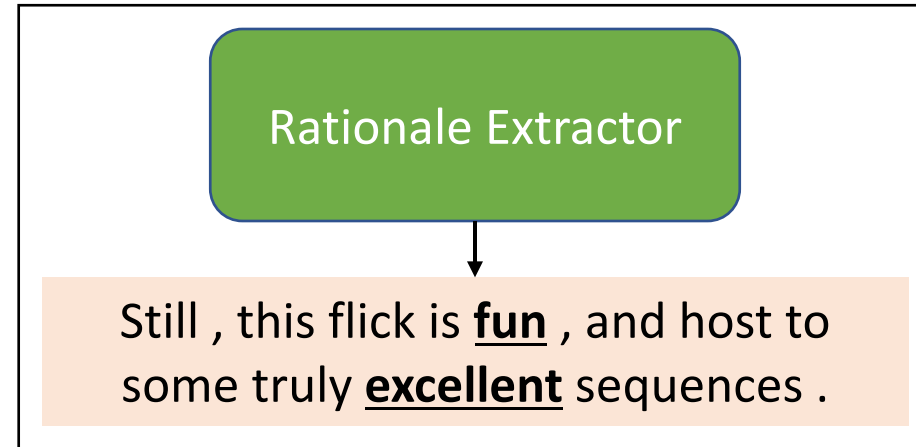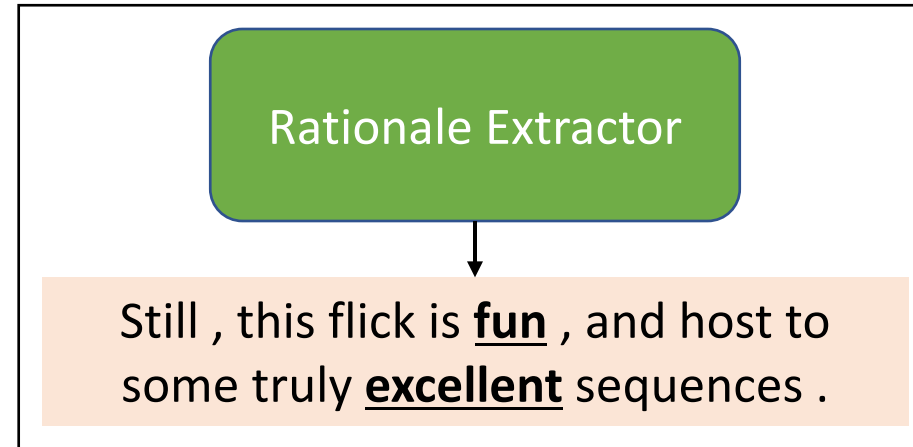
The explanation makes sense to us! 👍

# Three Desiderata of Rationale Extraction

**1. Faithfulness**

The explanation accurately reflects my reasoning process! 👍

Neural Language Model (NLM)

Rationale Extractor

Still , this flick is **fun** , and host to some truly **excellent** sequences .

**2. Plausibility**

The explanation makes sense to us! 👍

**3. Task Performance**

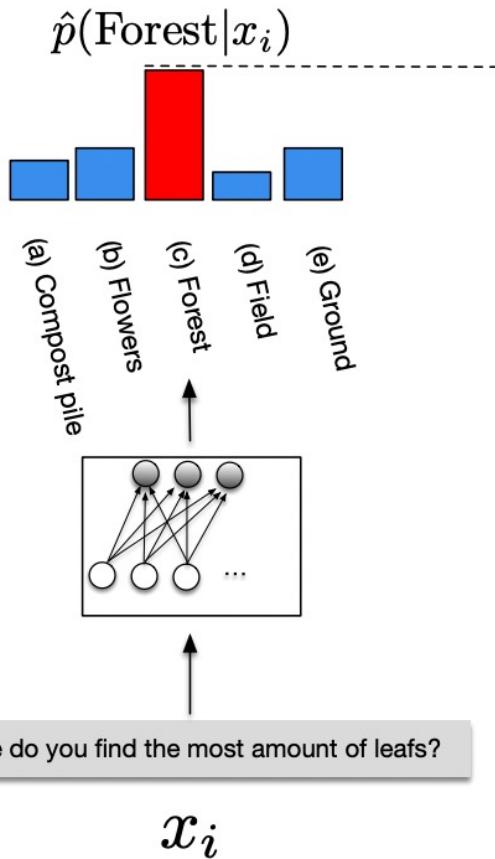| Score | BoolQ | CB | COPA | MultiRC |
|-------|-------|-----------|-------|-----------|
| 91.0 | 92.3 | 96.9/98.0 | 99.2 | 89.2/65.2 |
| 90.9 | 92.0 | 95.9/97.6 | 98.2 | 88.4/63.0 |
| 90.6 | 91.0 | 98.6/99.2 | 97.4 | 88.6/63.2 |
| 90.4 | 91.4 | 95.8/97.6 | 98.0 | 88.3/63.0 |
| 90.3 | 90.4 | 95.7/97.6 | 98.4 | 88.2/63.7 |
| 89.8 | 89.0 | 95.8/98.9 | 100.0 | 81.8/51.9 |

**Faithfulness**                    **Plausibility**

# Faithfulness

Comprehensiveness (Comp)    Sufficiency (Suff)

# Plausibility

# Faithfulness

**Comprehensiveness (Comp)**     **Sufficiency (Suff)**
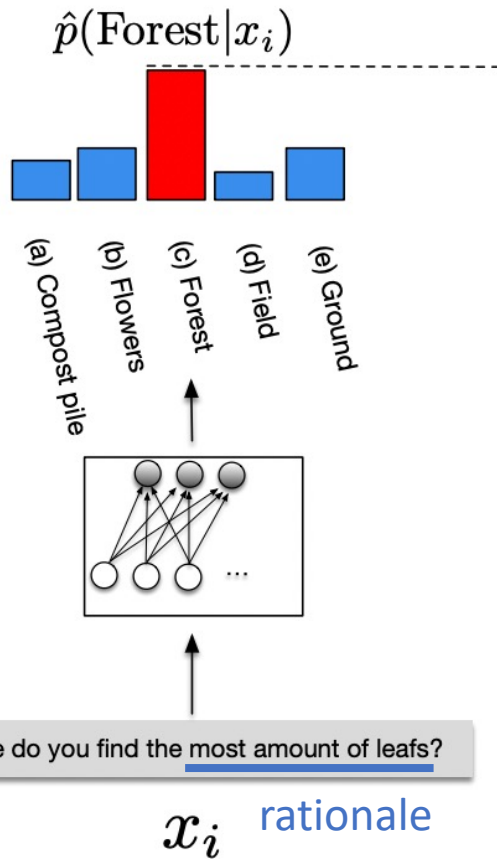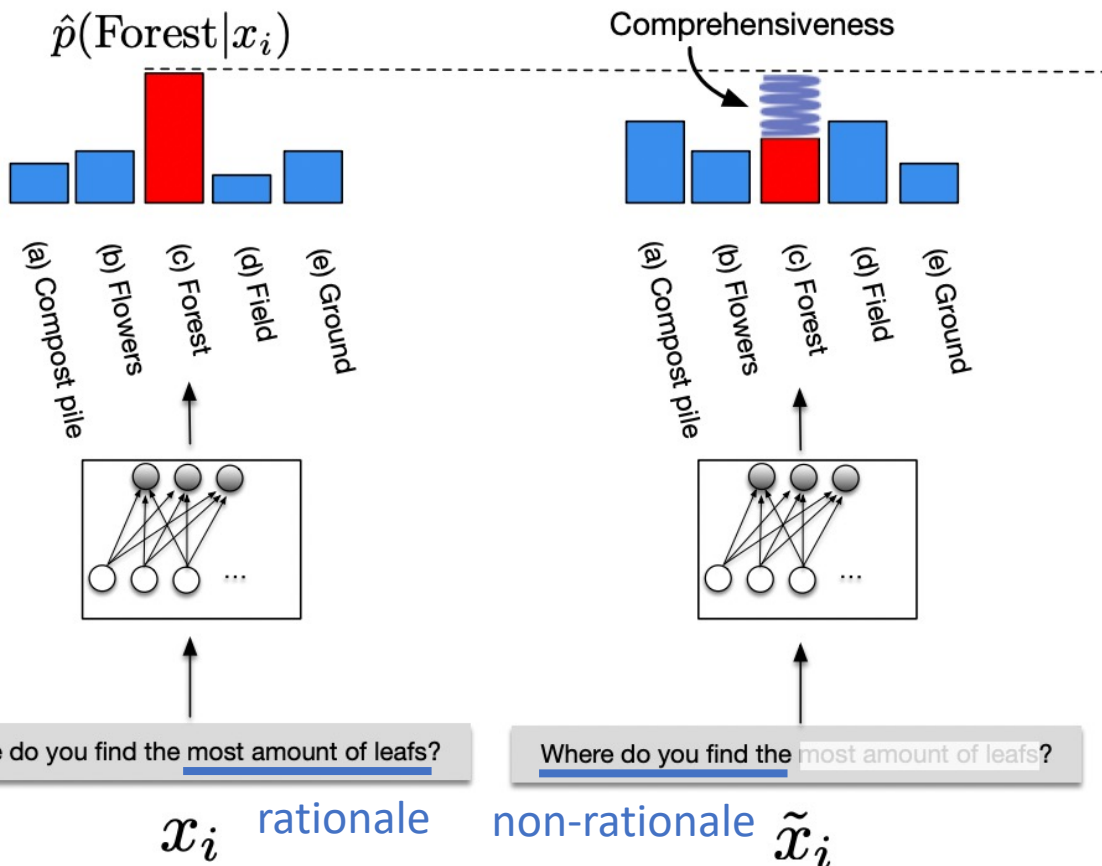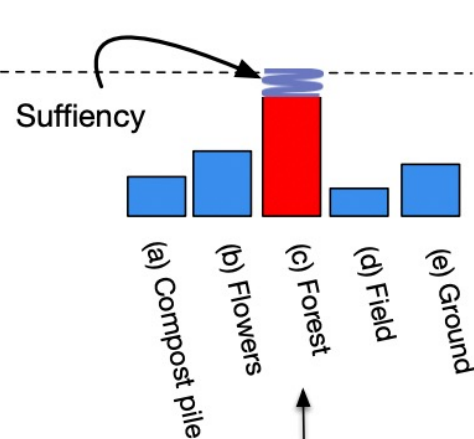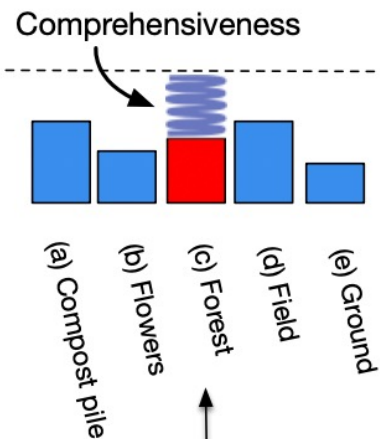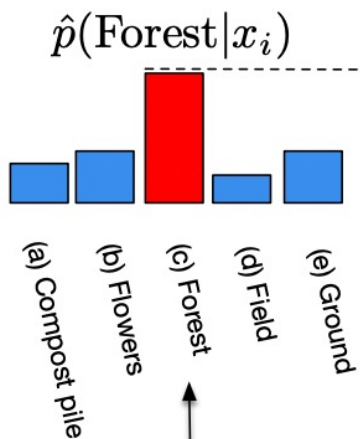
# Plausibility



$\hat{p}(\text{Forest}|x_i)$

(a) Compost pile
(b) Flowers
(c) Forest
(d) Field
(e) Ground

Where do you find the most amount of leafs?

$x_i$

# Faithfulness

**Comprehensiveness (Comp)**      **Sufficiency (Suff)**

# Plausibility

$\hat{p}(\text{Forest}|x_i)$

(a) Compost pile
(b) Flowers
(c) Forest
(d) Field
(e) Ground

Where do you find the most amount of leafs?

$x_i$    rationale

# Faithfulness

# Plausibility

## Comprehensiveness (Comp)     Sufficiency (Suff)



$\hat{p}(\text{Forest}|x_i)$

Comprehensiveness

(a) Compost pile
(b) Flowers
(c) Forest
(d) Field
(e) Ground

Where do you find the most amount of leafs?

rationale     non-rationale

$x_i$     $\tilde{x}_i$

*Higher* Comp is better!

# Faithfulness

## Comprehensiveness (Comp)    ## Sufficiency (Suff)

# Plausibility

$\hat{p}(\text{Forest}|x_i)$

Comprehensiveness

Suffiency

(a) Compost pile
(b) Flowers
(c) Forest
(d) Field
(e) Ground

Where do you find the most amount of leafs?

$x_i$    rationale

Where do you find the most amount of leafs?

$\tilde{x}_i$

Where do you find the most amount of leafs?
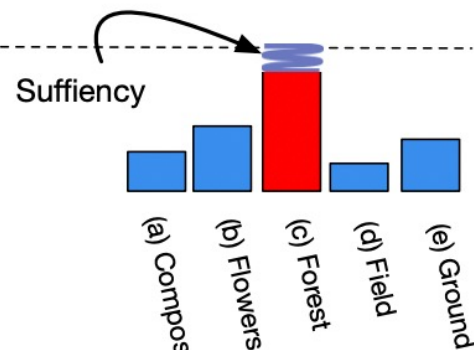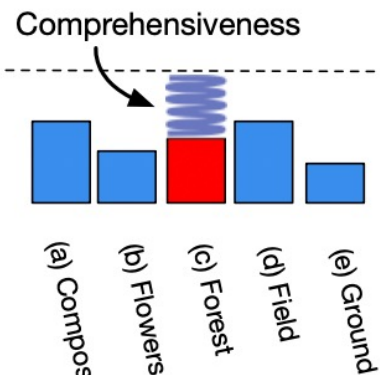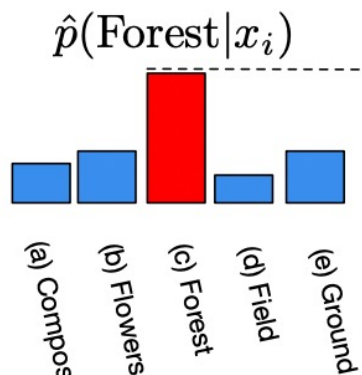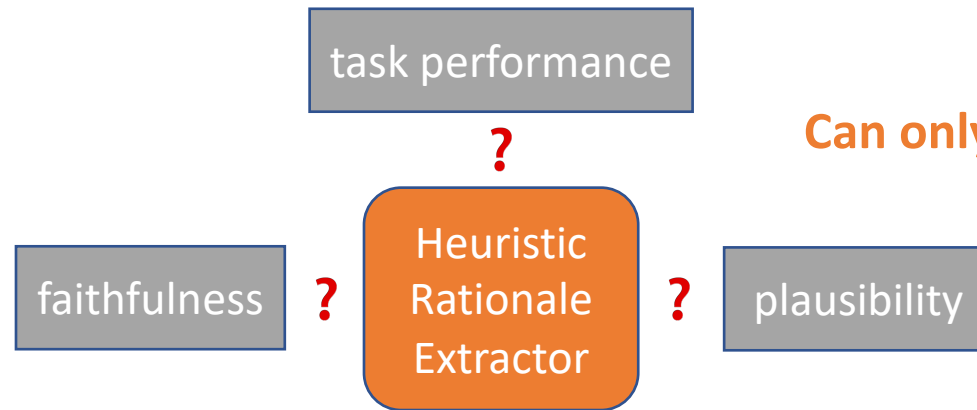
$r_i$    rationale

*Lower* Suff is better!

# Faithfulness

### Comprehensiveness (Comp)     ### Sufficiency (Suff)

# Plausibility

### Similarity to Gold Rationales

$\hat{p}(\text{Forest}|x_i)$

Comprehensiveness

Suffiency

(a) Compost pile
(b) Flowers
(c) Forest
(d) Field
(e) Ground

Where do you find the most amount of leafs?

$x_i$

$\tilde{x}_i$

$r_i$

Predicted Rationale

**Still** , this **flick** is **fun** , and **host** to **some** truly excellent sequences .

Gold Rationale

Still , this **flick** is **fun** , and host to some **truly excellent sequences** .

# UNIREX: UNIfied Learning Framework for Rationale EXtraction

**Existing Works**



task performance

**?**

faithfulness   **?**   Heuristic Rationale Extractor   **?**   plausibility

**Can only choose up to two!**

# UNIREX: UNIfied Learning Framework for Rationale EXtraction

**Existing Works**

task performance

**?**

faithfulness **?** Heuristic Rationale Extractor **?** plausibility

**Can only choose up to two!**

**UNIREX**

task performance

**Can optimize for all three!**
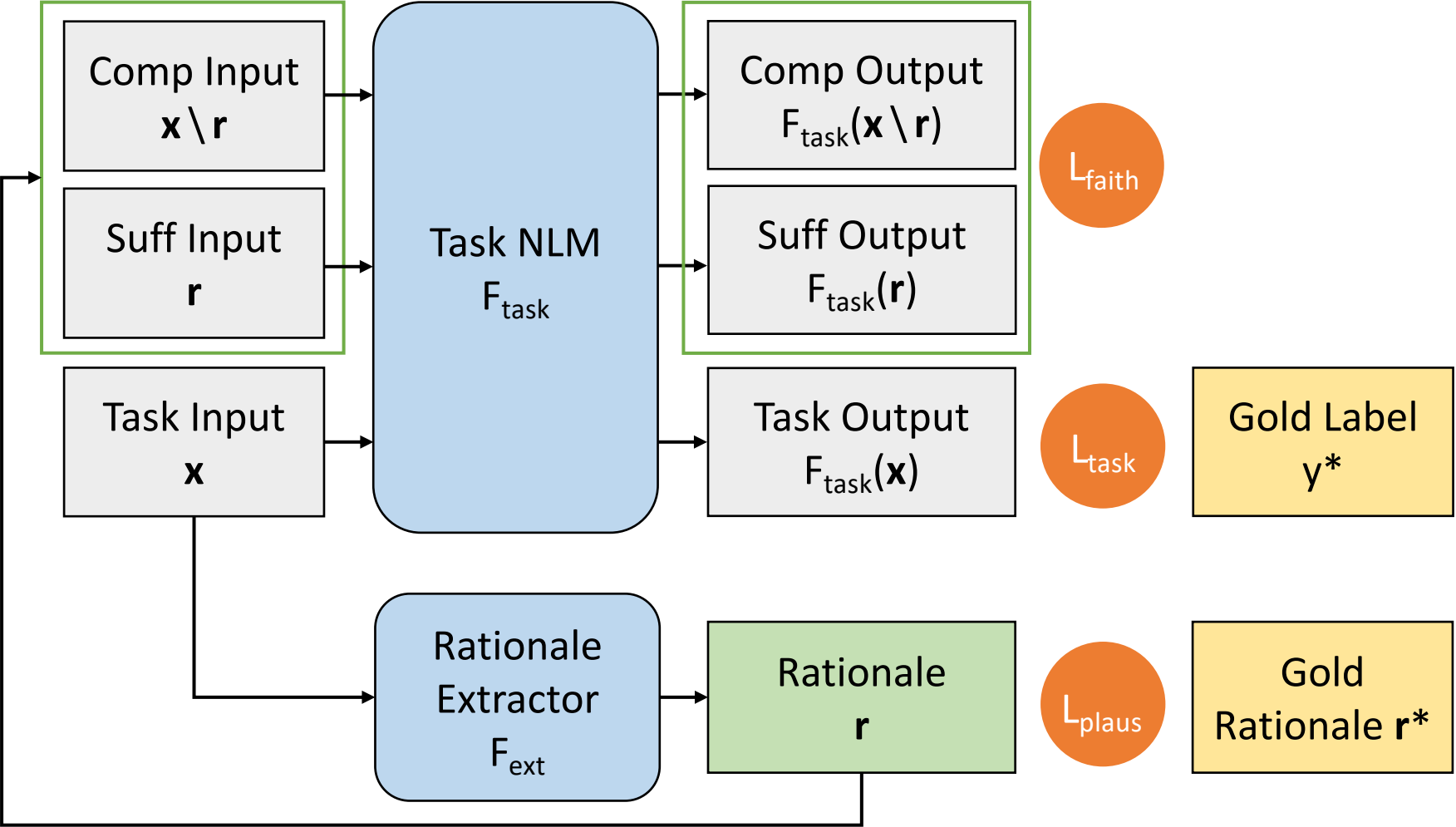
faithfulness Learned Rationale Extractor plausibility
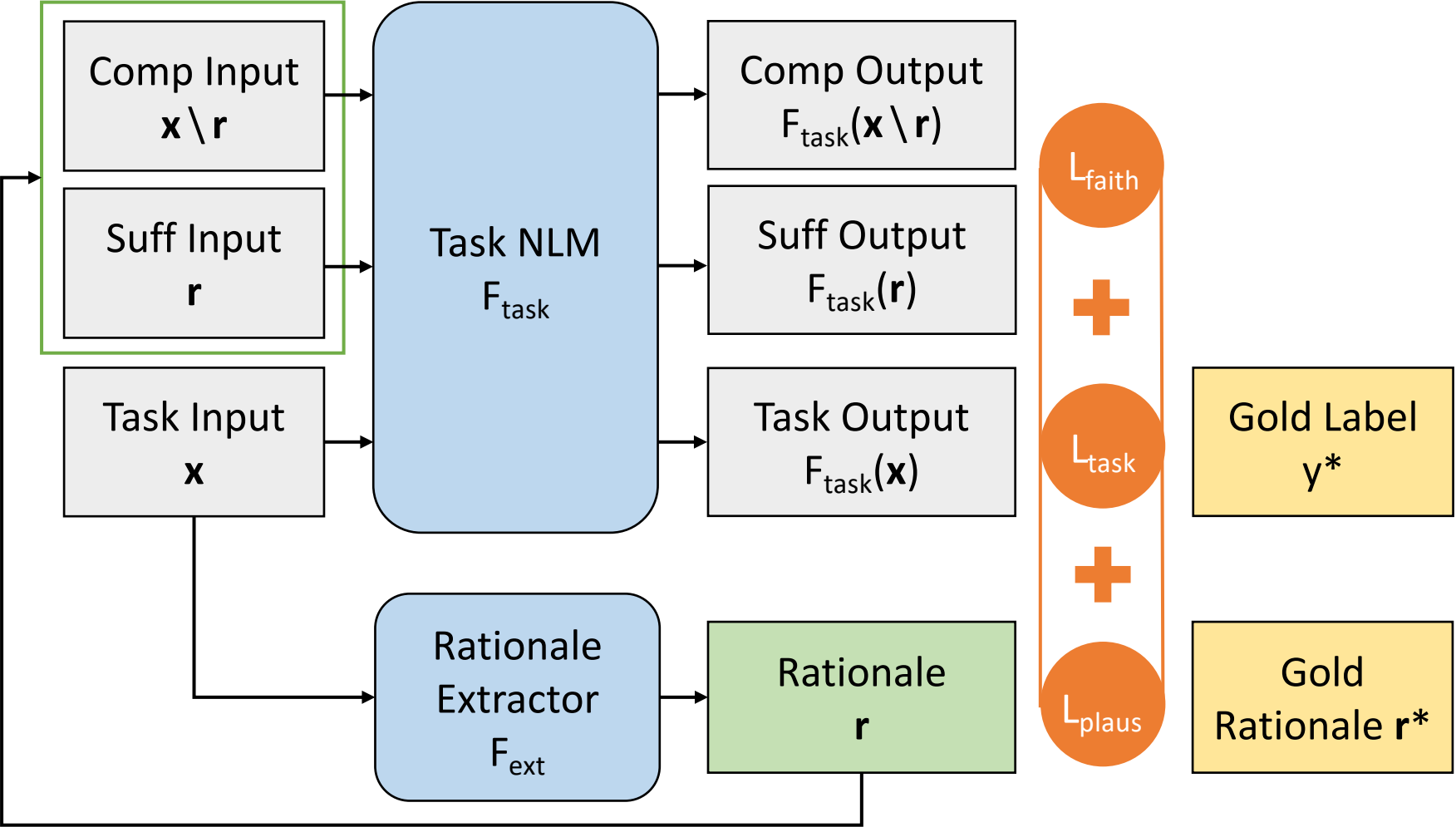
# UNIREX

# UNIREX

# UNIREX

# UNIREX

# UNIREX

# Experiments: Main Results
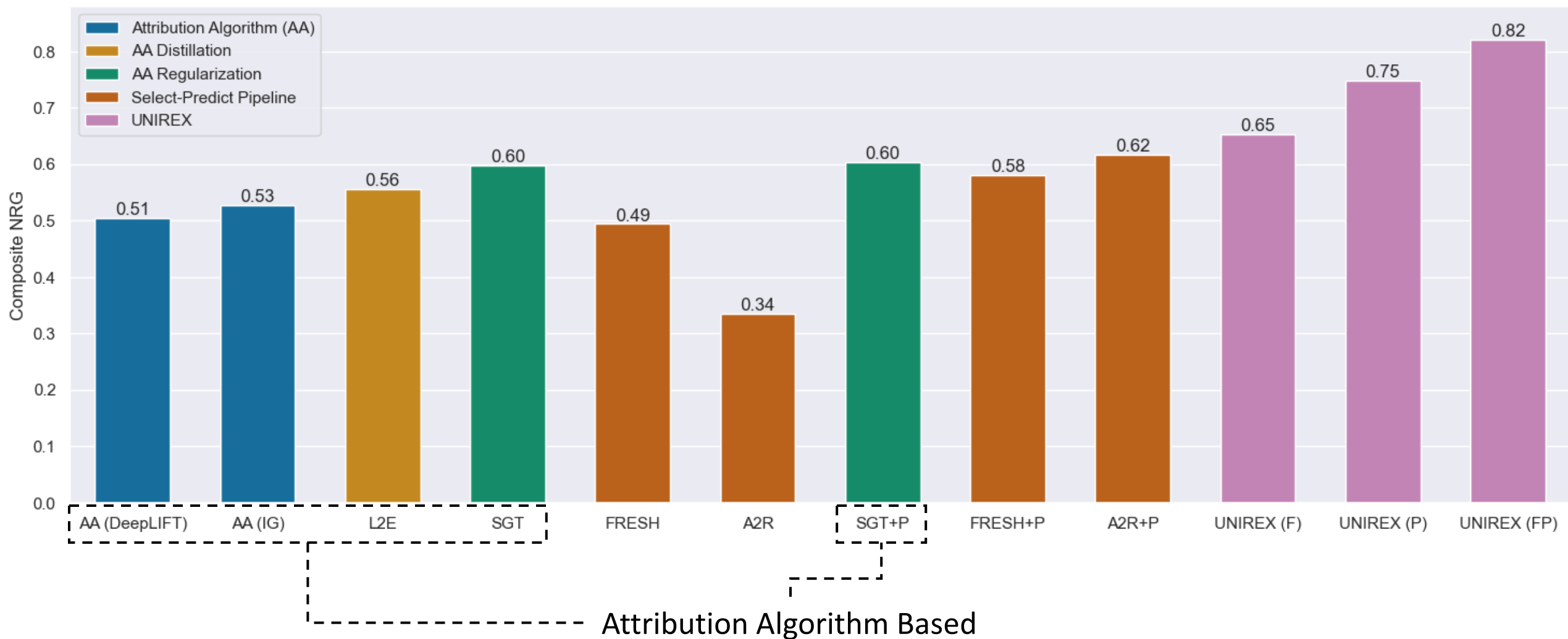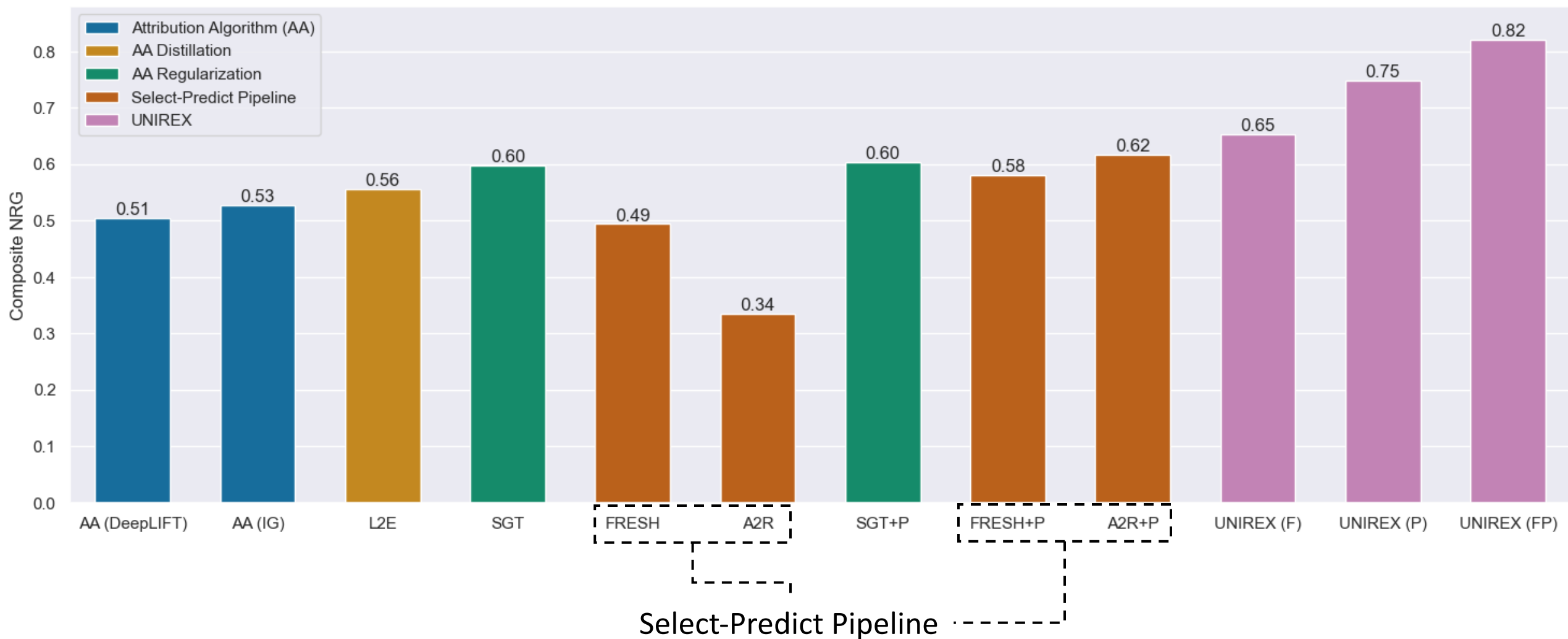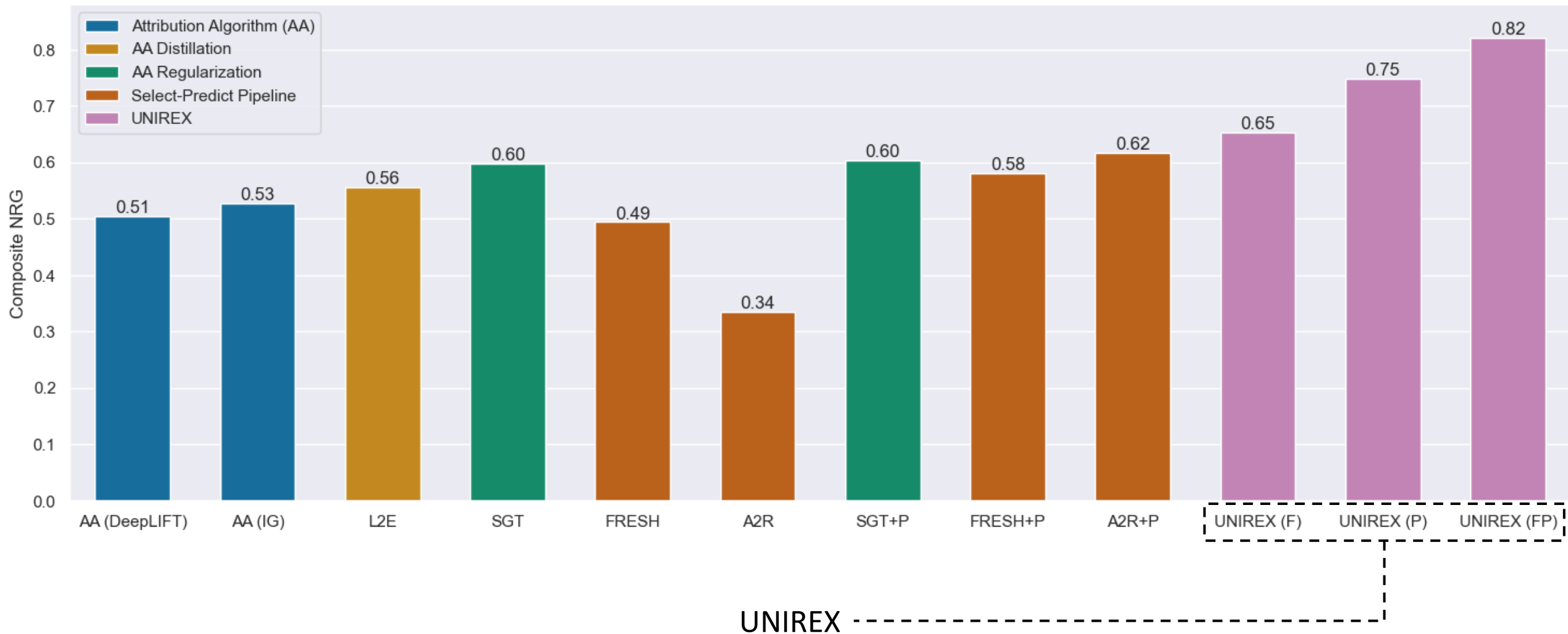
# Experiments: Main Results

# Experiments: Main Results

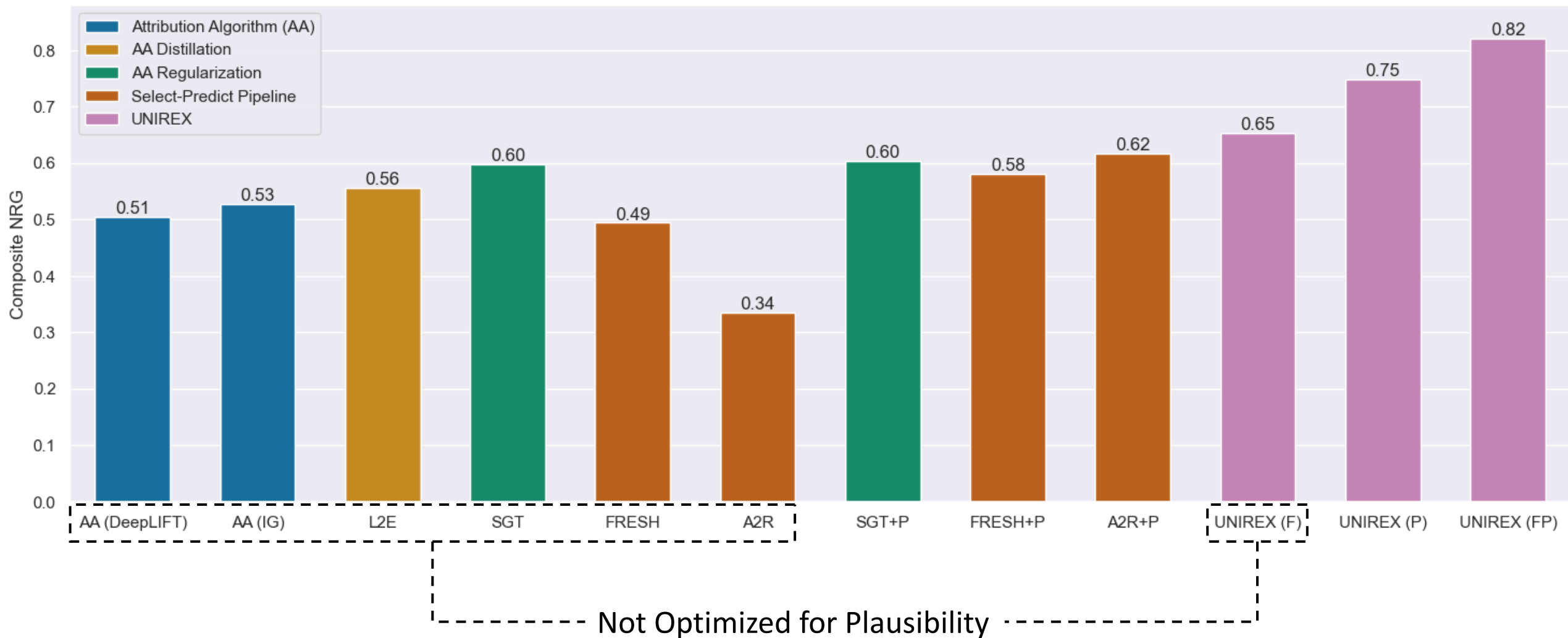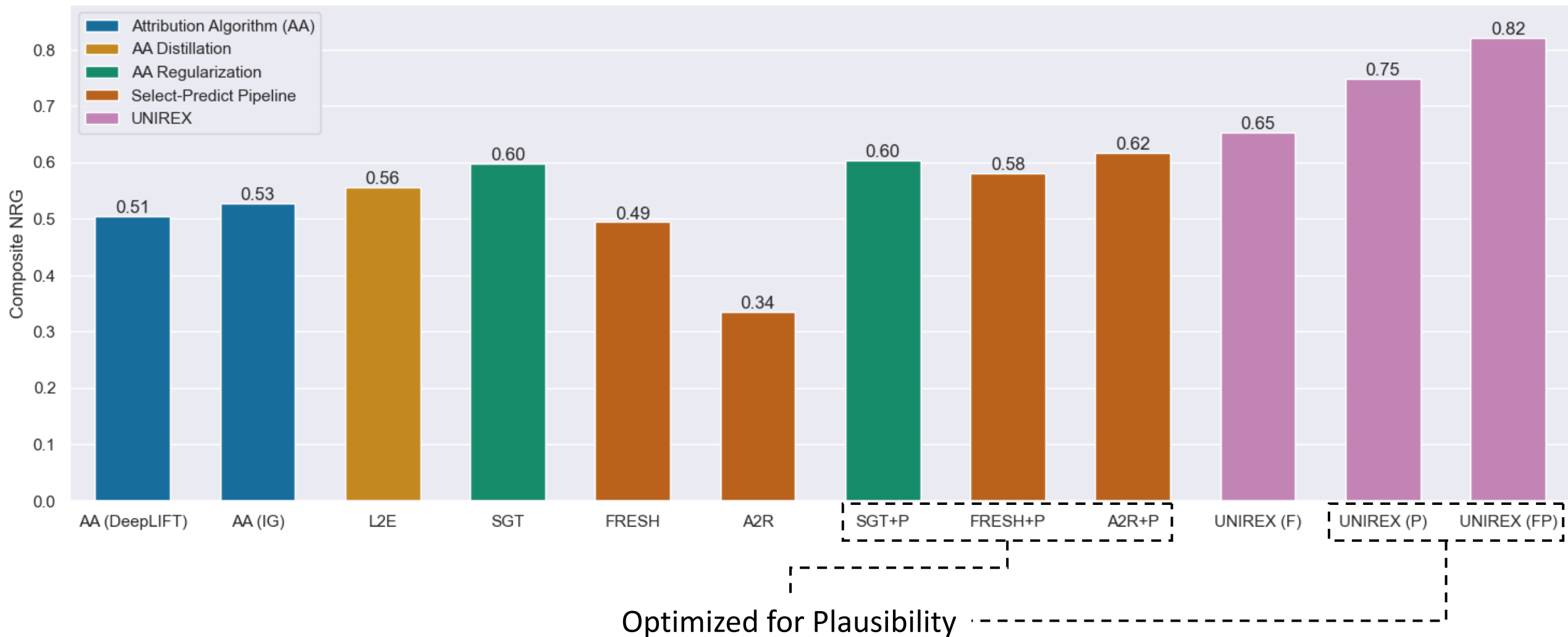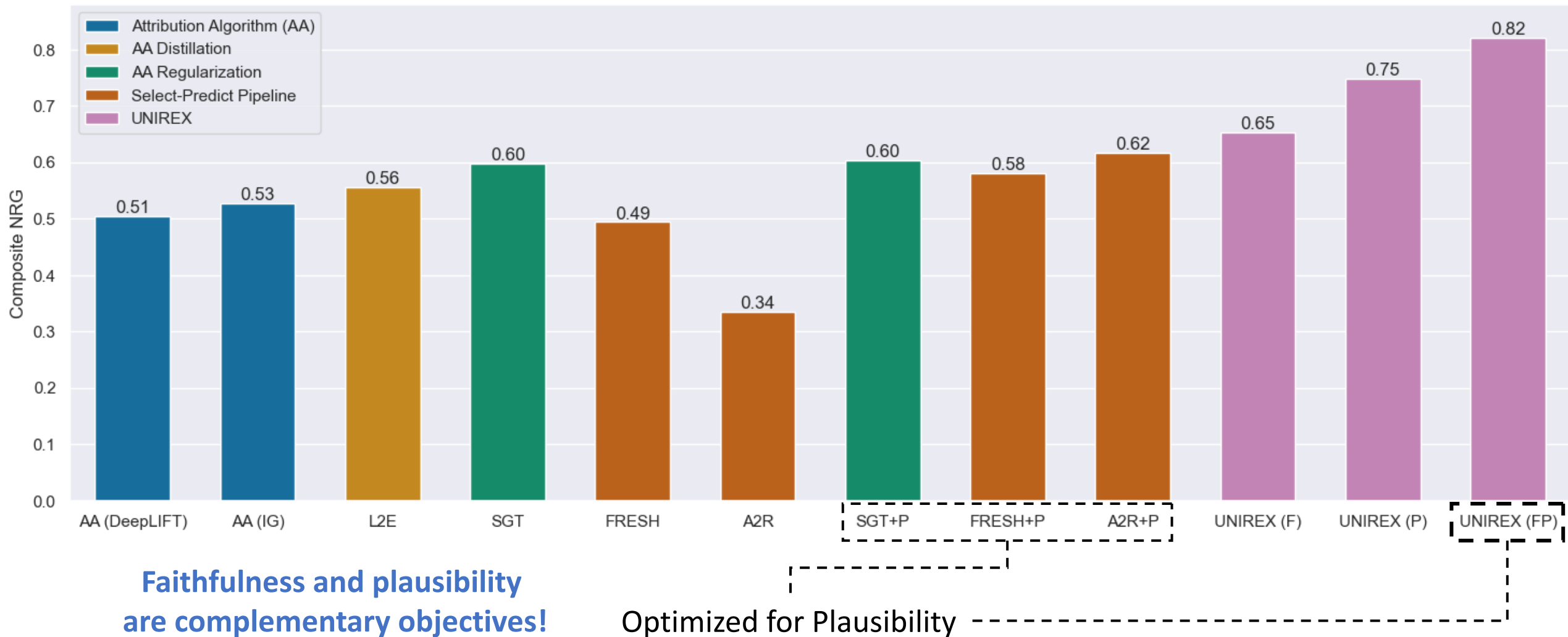# Experiments: Main Results

# Experiments: Main Results

# Experiments: Main Results

Optimized for Plausibility

# Experiments: Main Results

**Faithfulness and plausibility
are complementary objectives!**

Optimized for Plausibility

# Thank You!