

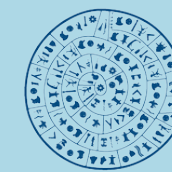


# Highways of Horizontal Gene Transfer for Modeling Large Gene Flow Events Under the UndatedDTL Model

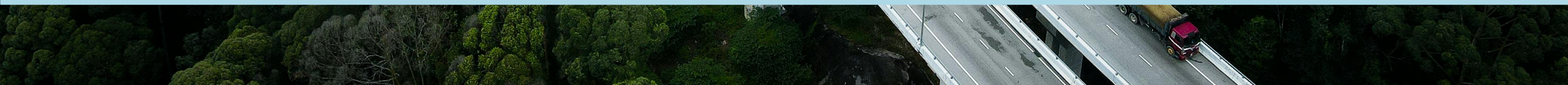
**Noah A. Wahl**, Benoit Morel, Oliver Schick, Alexandros Stamatakis,  
Tom A. Williams, Gergely J. Szöllősi (TBF)



ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΡΗΤΗΣ  
UNIVERSITY OF CRETE



**FORTH**  
INSTITUTE OF COMPUTER SCIENCE

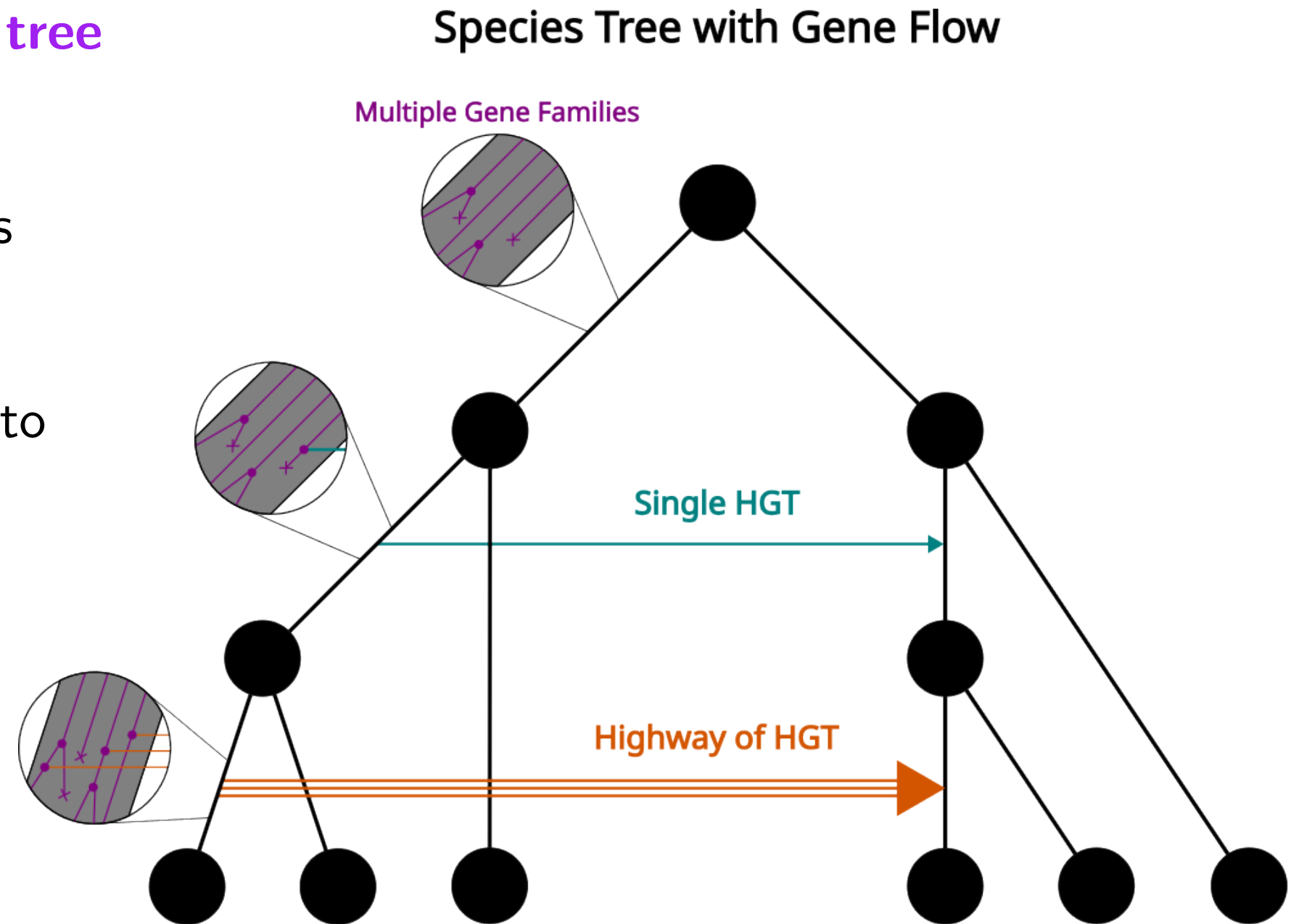


# Gene Tree Reconciliation and Highways

- Species and gene tree topologies often disagree
- Given a species tree and a collection of **gene tree samples** of gene families
- We try to infer a series of gene origination, **duplication, loss and transfer** (DTL) events under maximum likelihood
- The **UndatedDTL** (multi-)model uses branch-wise DTL and origination parameters to sample these events

# Gene Tree Reconciliation and Highways

- Species and gene tree topologies often disagree
- Given a species tree and a collection of **gene tree samples** of gene families
- We try to infer a series of gene origination, **duplication, loss and transfer** (DTL) events under maximum likelihood
- The **UndatedDTL** (multi-)model uses branch-wise DTL and origination parameters to sample these events
- Highways of HGT are pairs of species tree branches that have an **elevated** rate of gene exchange [1]
- Often an indicator of symbiont genes being transferred to the host genome



\*HGT = Horizontal Gene Transfers

[1] Beiko, Robert G., Timothy J. Harlow, and Mark A. Ragan. "Highways of gene sharing in prokaryotes." Proceedings of the National Academy of Sciences 102.40 (2005)

# Motivation

- **Horizontal gene transfer** (HGT) makes a lot of things more challenging
- The transfer rate needs to consider **all possible transfer targets**
- Global or family-wise DTL rates can be inferred even for large trees
- Branch-wise rates for small trees
- Ideal case: **donor** and **recipient** rates for all species tree branches
- Instead: can we just identify **highways** to deal with the most extreme outliers?

# Motivation

- **Horizontal gene transfer** (HGT) makes a lot of things more challenging
- The transfer rate needs to consider **all possible transfer targets**
- Global or family-wise DTL rates can be inferred even for large trees
- Branch-wise rates for small trees
- Ideal case: **donor** and **recipient** rates for all species tree branches
- Instead: can we just identify **highways** to deal with the most extreme outliers?
  
- Extract that pair as a separate parameter  $H_{e,d}$

# Extending the UndatedDTL Model in AleRax [2] with Highways

Collection of gene families with a distribution of gene trees per family  
1 species tree  
D, T, L, **H** branch-wise parameters



Gene tree distributions get turned into clade splits and conditional clade probabilities (CCPs) [3,4]



Double recursion over species tree branches ( $e$ ) and clade splits ( $\gamma', \gamma'' \mid \gamma$ ) in the gene trees to compute conditional likelihoods  $\Pi(e, \gamma)$

[2] Morel, Benoit, et al. "AleRax: a tool for gene and species tree co-estimation and reconciliation under a probabilistic model of gene duplication, transfer, and loss." *Bioinformatics* 40.4 (2024)

[3] Höhna, Sebastian, and Alexei J. Drummond. "Guided tree topology proposals for Bayesian phylogenetic inference." *Systematic biology* 61.1 (2012)

[4] Larget, Bret. "The estimation of tree posterior probabilities using conditional clade probability distributions." *Systematic biology* 62.4 (2013)

# Extending the UndatedDTL Model in AleRax [2] with Highways

Collection of gene families with a distribution of gene trees per family  
1 species tree  
D, T, L, **H** branch-wise parameters



Gene tree distributions get turned into clade splits and conditional clade probabilities (CCPs) [3,4]



Double recursion over species tree branches ( $e$ ) and clade splits ( $\gamma', \gamma'' \mid \gamma$ ) in the gene trees to compute conditional likelihoods  $\Pi(e, \gamma)$

UndatedDTL:

$$\begin{aligned}\Pi(e, \gamma) &= S + SL \\ &+ D + DL \\ &+ T + TL\end{aligned}$$

UndatedDTL + **H**:

$$\begin{aligned}\Pi(e, \gamma) &= S + SL \\ &+ D + DL \\ &+ T + TL \\ &+ H + HL\end{aligned}$$

[2] Morel, Benoit, et al. "AleRax: a tool for gene and species tree co-estimation and reconciliation under a probabilistic model of gene duplication, transfer, and loss." *Bioinformatics* 40.4 (2024)

[3] Höhna, Sebastian, and Alexei J. Drummond. "Guided tree topology proposals for Bayesian phylogenetic inference." *Systematic biology* 61.1 (2012)

[4] Larget, Bret. "The estimation of tree posterior probabilities using conditional clade probability distributions." *Systematic biology* 62.4 (2013)

# Extending the UndatedDTL Model in AleRax [2] with Highways

Collection of gene families with a distribution of gene trees per family  
1 species tree  
D, T, L, **H** branch-wise parameters



Gene tree distributions get turned into clade splits and conditional clade probabilities (CCPs) [3,4]



Double recursion over species tree branches ( $e$ ) and clade splits ( $\gamma', \gamma'' \mid \gamma$ ) in the gene trees to compute conditional likelihoods  $\Pi(e, \gamma)$

UndatedDTL:

$$\begin{aligned}\Pi(e, \gamma) = & S + SL \\ & + D + DL \\ & + T + TL\end{aligned}$$

UndatedDTL + **H**:

$$\begin{aligned}\Pi(e, \gamma) = & S + SL \\ & + D + DL \\ & + T + TL \\ & + H + HL\end{aligned}$$



Easy, right?

[2] Morel, Benoit, et al. "AleRax: a tool for gene and species tree co-estimation and reconciliation under a probabilistic model of gene duplication, transfer, and loss." *Bioinformatics* 40.4 (2024)

[3] Höhna, Sebastian, and Alexei J. Drummond. "Guided tree topology proposals for Bayesian phylogenetic inference." *Systematic biology* 61.1 (2012)

[4] Larget, Bret. "The estimation of tree posterior probabilities using conditional clade probability distributions." *Systematic biology* 62.4 (2013)

# Overfitting

- The UndatedDTL model can already have many parameters (branch-wise, per-family)
- But those are **user-regulated**
- Each highway introduces 1 additional parameter
- Highway parameters are chosen **during optimization**
- Too many (unnecessary) parameters slows down the optimizer significantly
- Need good heuristics with penalties for adding too many parameters and avoid **overfitting**

# Overfitting

- The UndatedDTL model can already have many parameters (branch-wise, per-family)
- But those are **user-regulated**
- Each highway introduces 1 additional parameter
- Highway parameters are chosen **during optimization**
- Too many (unnecessary) parameters slows down the optimizer significantly
- Need good heuristics with penalties for adding too many parameters and avoid **overfitting**
  
- Easy case: user provides highway candidates to be **validated**
- Have at most this many highway parameters

# Overfitting

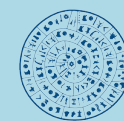
- The UndatedDTL model can already have many parameters (branch-wise, per-family)
- But those are **user-regulated**
- Each highway introduces 1 additional parameter
- Highway parameters are chosen **during optimization**
- Too many (unnecessary) parameters slows down the optimizer significantly
- Need good heuristics with penalties for adding too many parameters and avoid **overfitting**
  
- Easy case: user provides highway candidates to be **validated**
- Have at most this many highway parameters

This requires a lot of domain knowledge and good hypotheses!

# Highway Inference Heuristics

A good highway inference heuristic needs to achieve 2 things:

1. Find candidate node pairs with an exceptional number of transfers between them
2. Limit and test candidates quickly and discard ones with insufficient likelihood improvements



# Highway Inference Heuristics

A good highway inference heuristic needs to achieve 2 things:

1. Find candidate node pairs with an exceptional number of transfers between them
2. Limit and test candidates quickly and discard ones with insufficient likelihood improvements

Optimize DTL rates

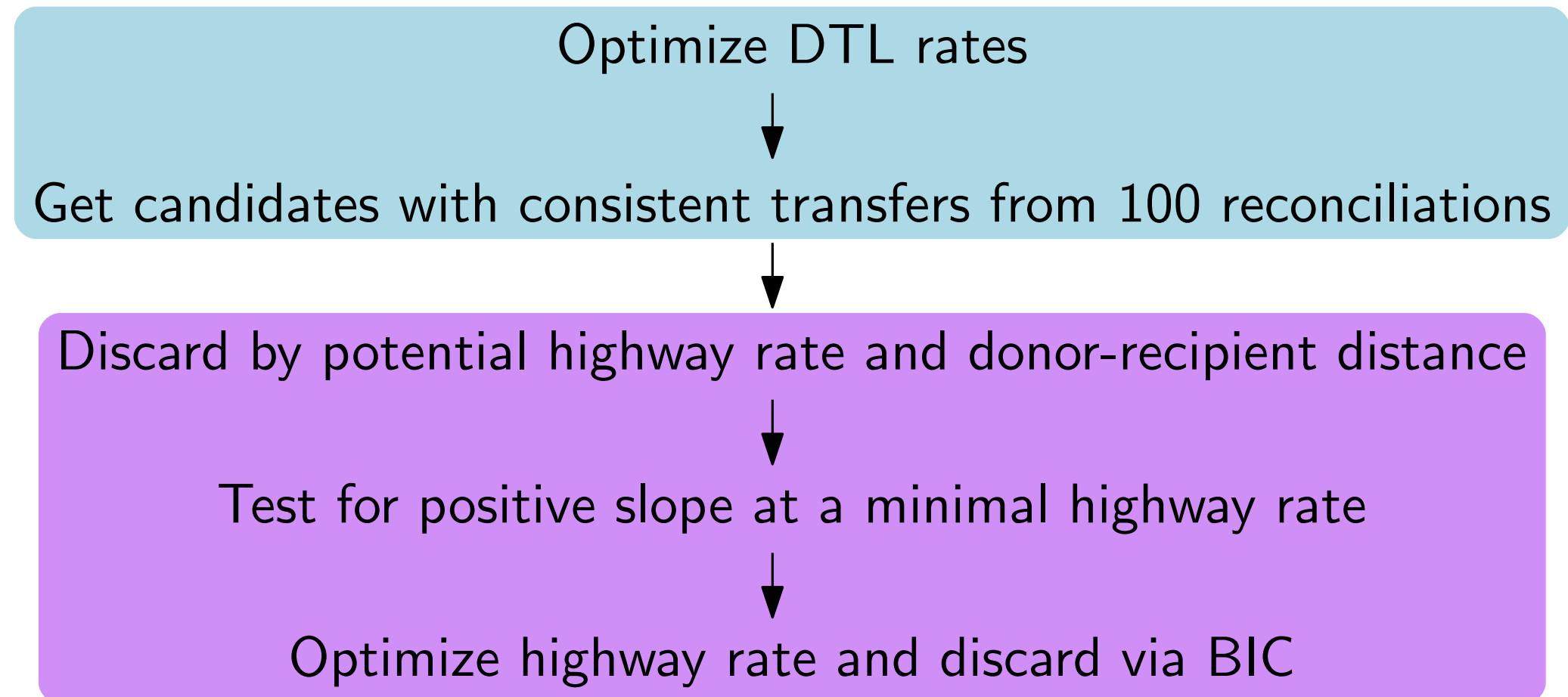


Get candidates with consistent transfers from 100 reconciliations

# Highway Inference Heuristics

A good highway inference heuristic needs to achieve 2 things:

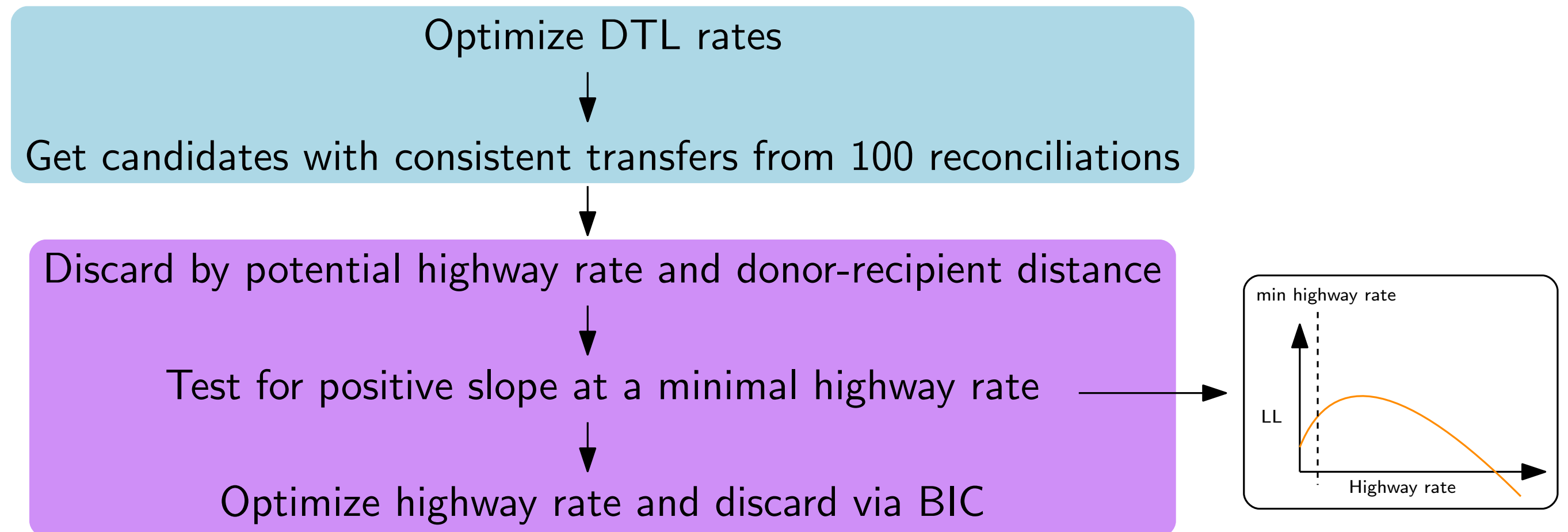
1. Find candidate node pairs with an exceptional number of transfers between them
2. Limit and test candidates quickly and discard ones with insufficient likelihood improvements



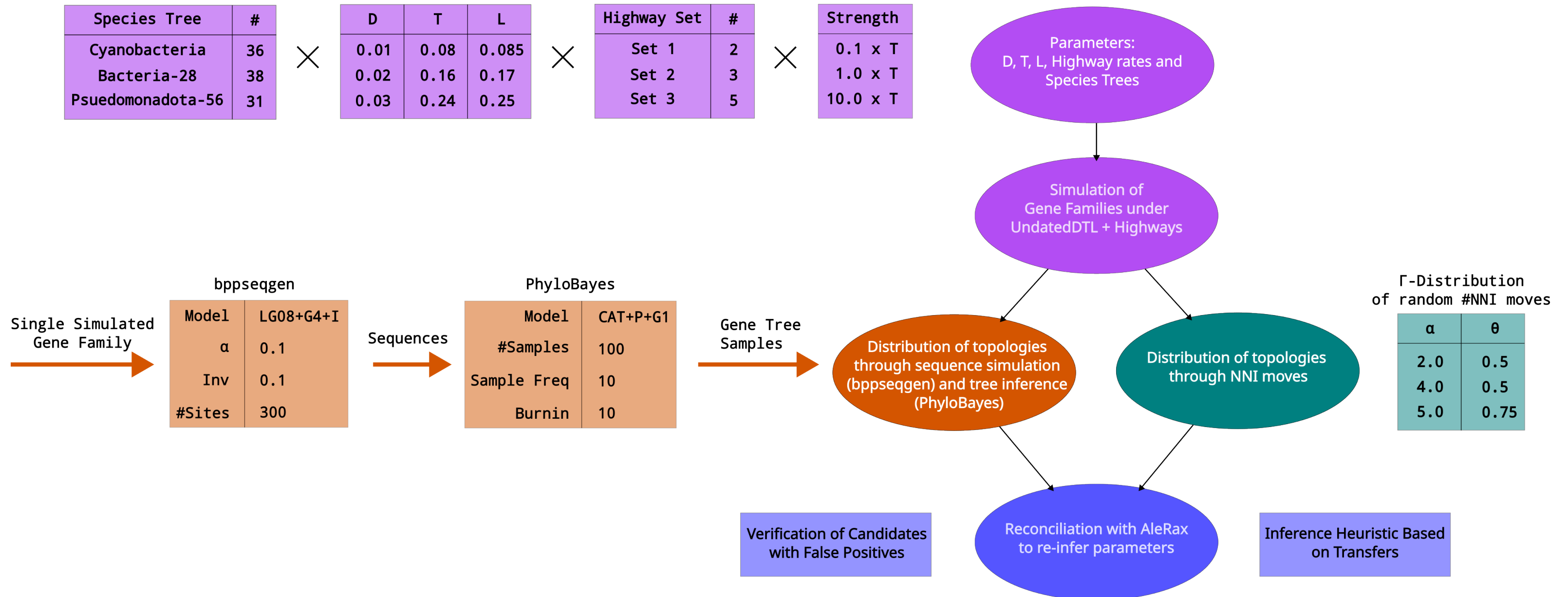
# Highway Inference Heuristics

A good highway inference heuristic needs to achieve 2 things:

1. Find candidate node pairs with an exceptional number of transfers between them
2. Limit and test candidates quickly and discard ones with insufficient likelihood improvements

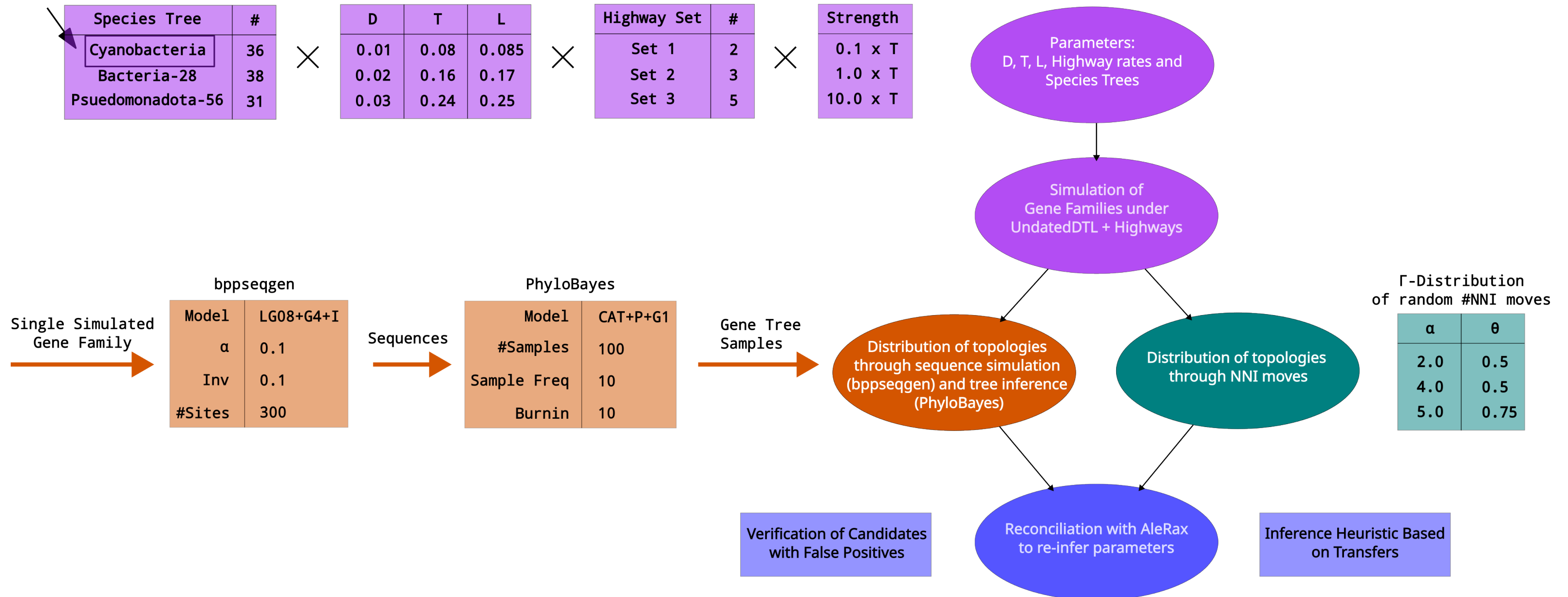


# Simulation Pipeline



# Simulation Pipeline

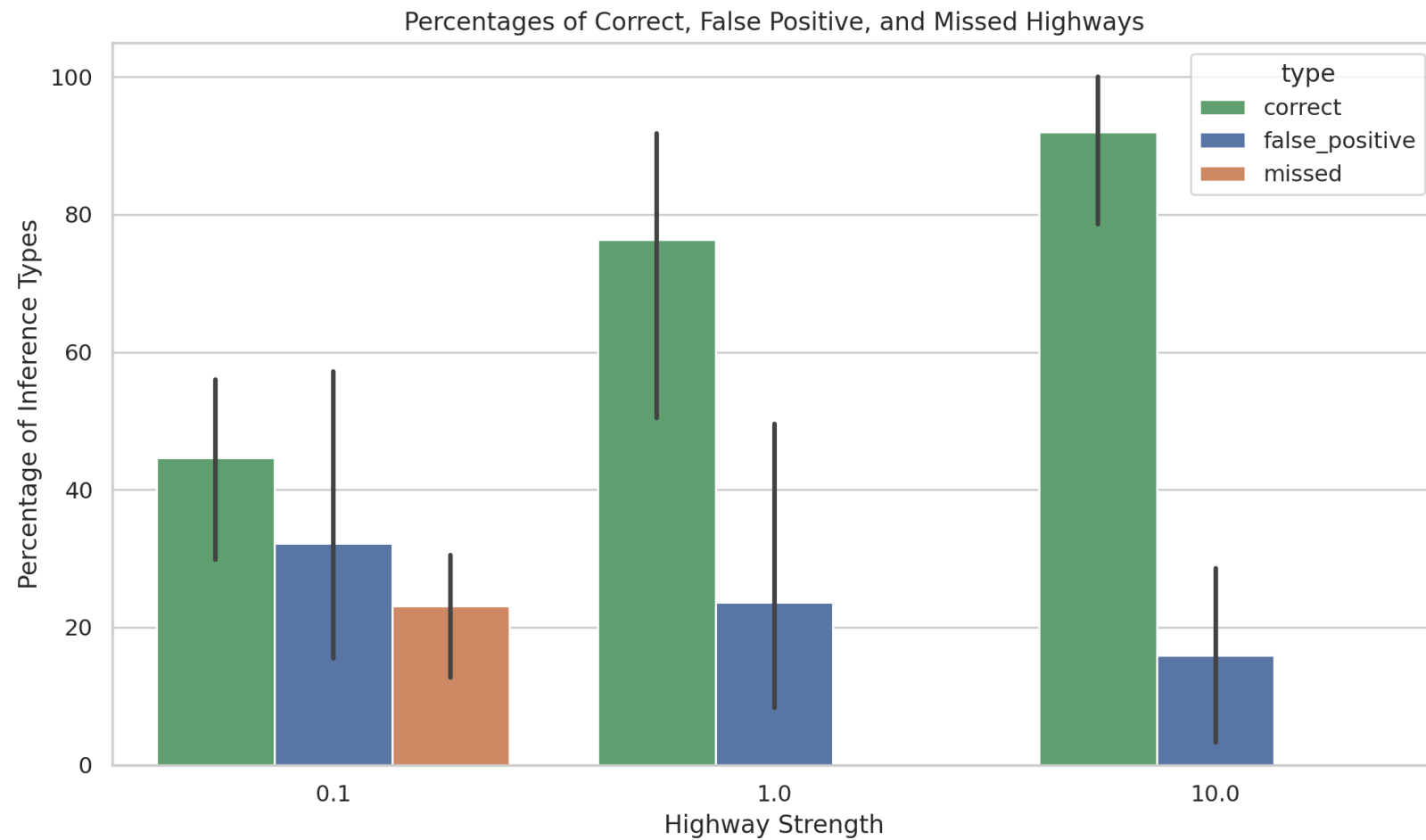
focus on these results



# Results (Transfer Heuristic)

$$\text{BIC: } \alpha \times \Delta LL > \log(N)$$

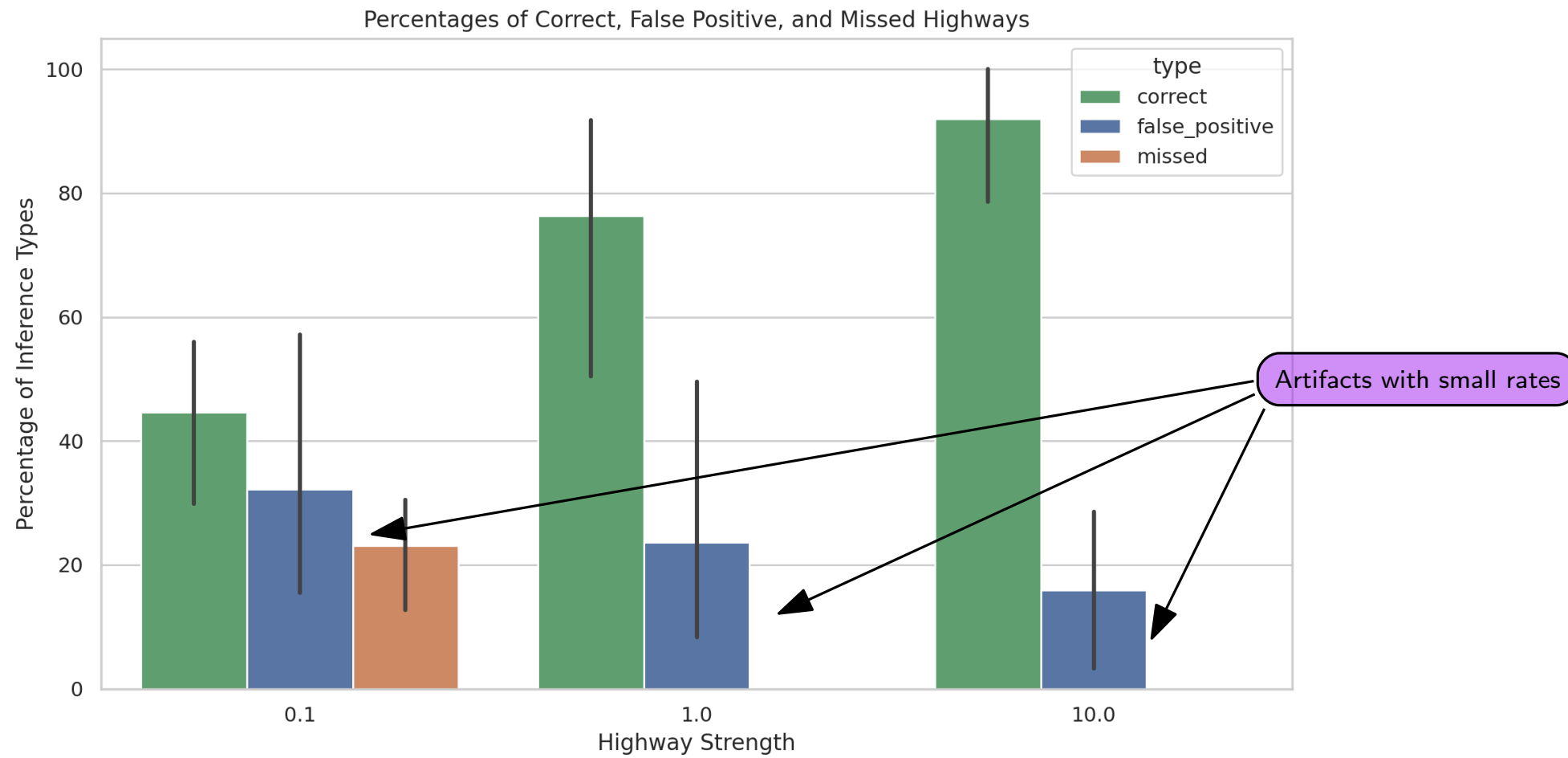
Weak pruning (distance, LL improvement, BIC  $\alpha = 2$ )



# Results (Transfer Heuristic)

$$\text{BIC: } \alpha \times \Delta LL > \log(N)$$

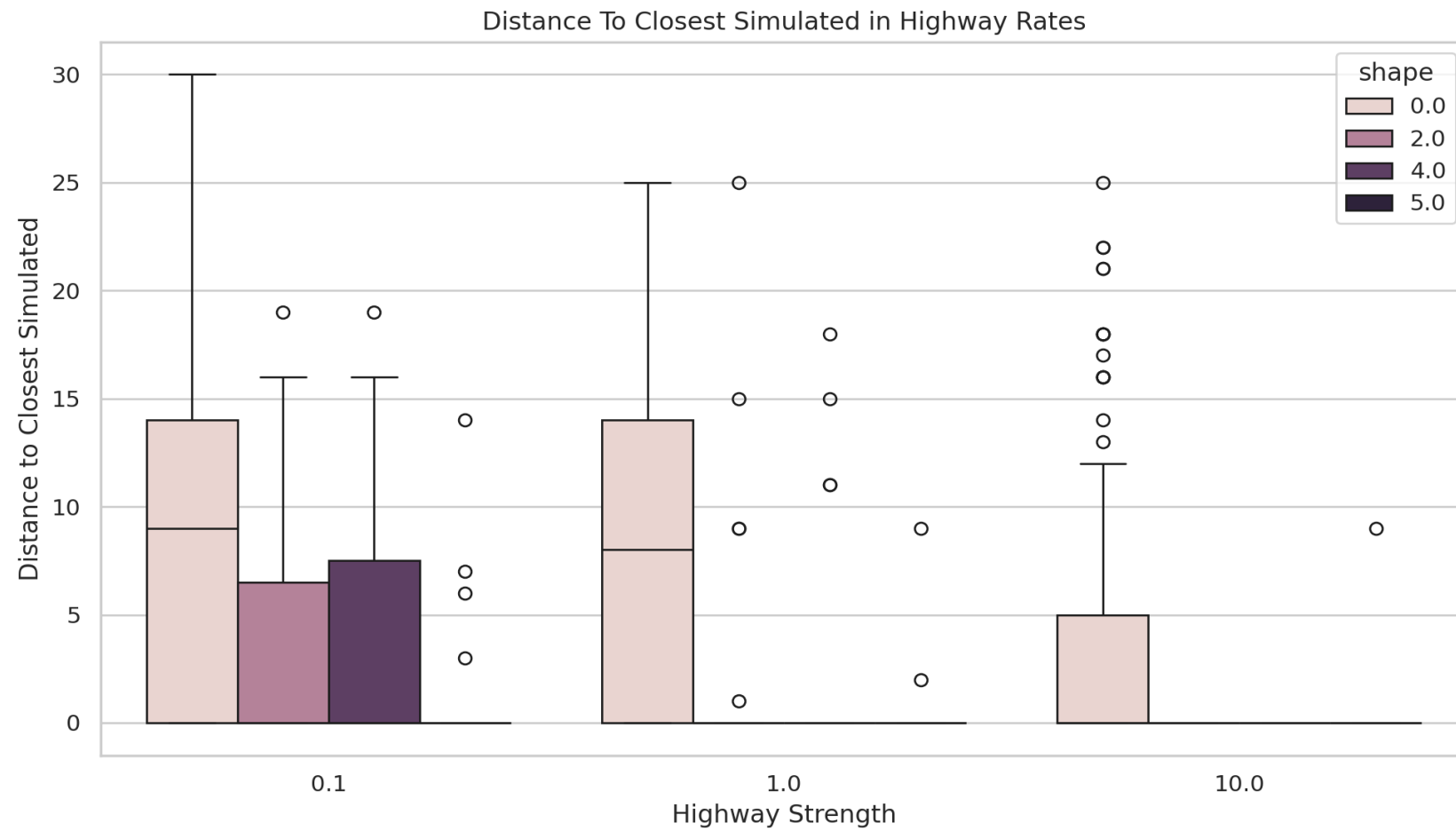
Weak pruning (distance, LL improvement, BIC  $\alpha = 2$ )



# Results (Transfer Heuristic)

$$\text{BIC: } \alpha \times \Delta LL > \log(N)$$

Weak pruning (distance, LL improvement, BIC  $\alpha = 2$ )

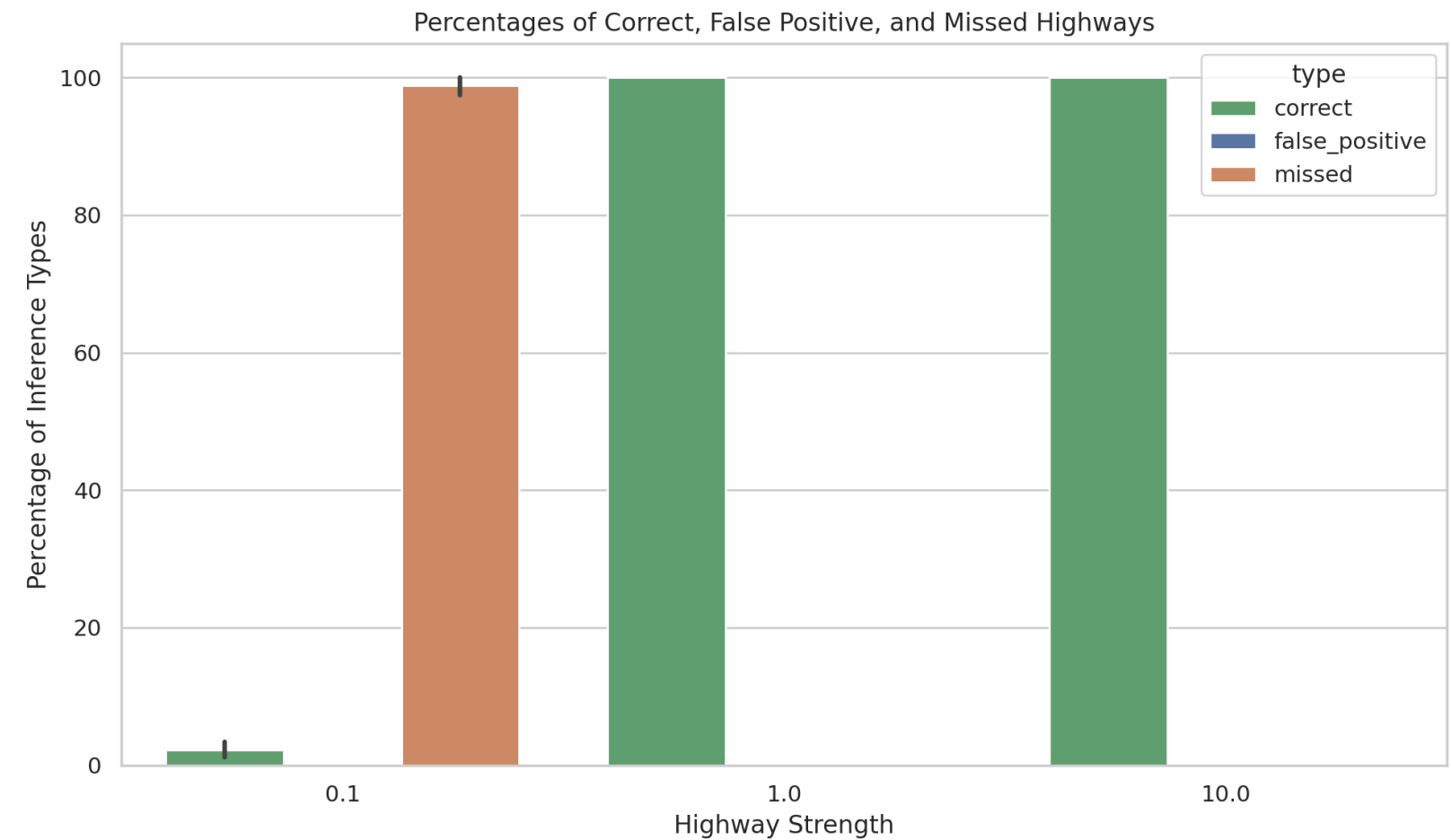
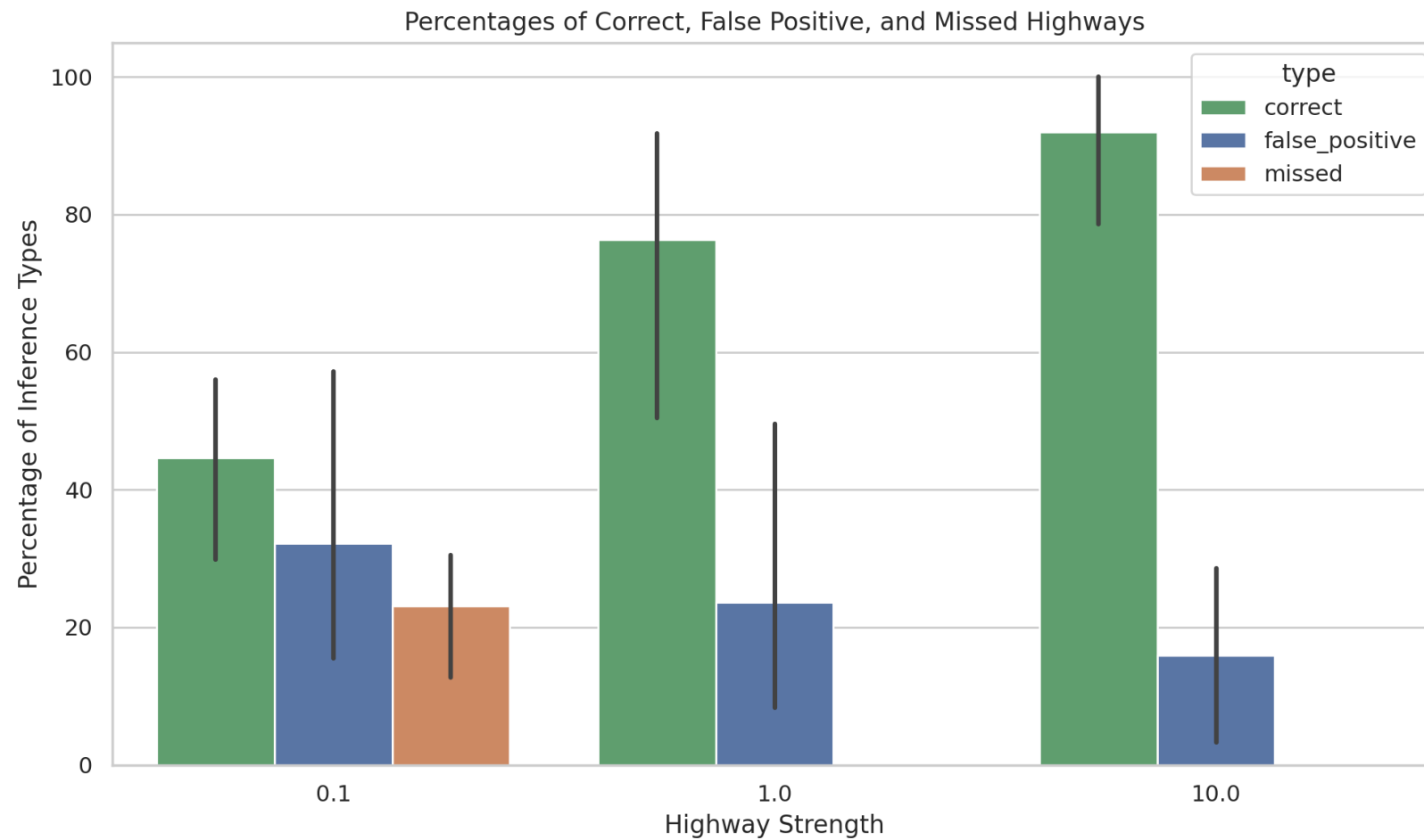


# Results (Transfer Heuristic)

$$\text{BIC: } \alpha \times \Delta LL > \log(N)$$

Weak pruning (distance, LL improvement, BIC  $\alpha = 2$ )

Strong pruning (distance, **potential highway rate**, LL improvement, BIC  $\alpha = 1$ )

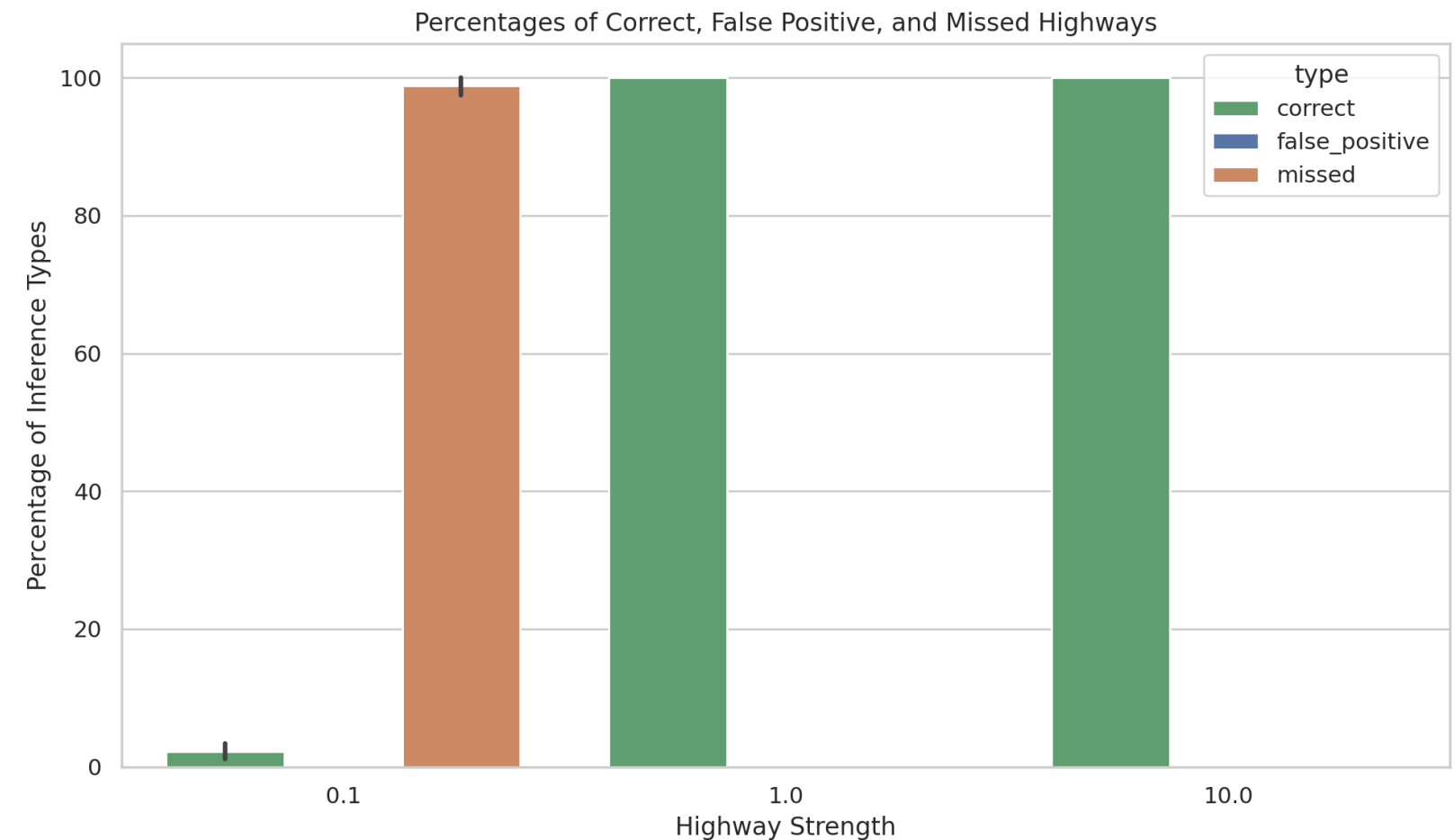
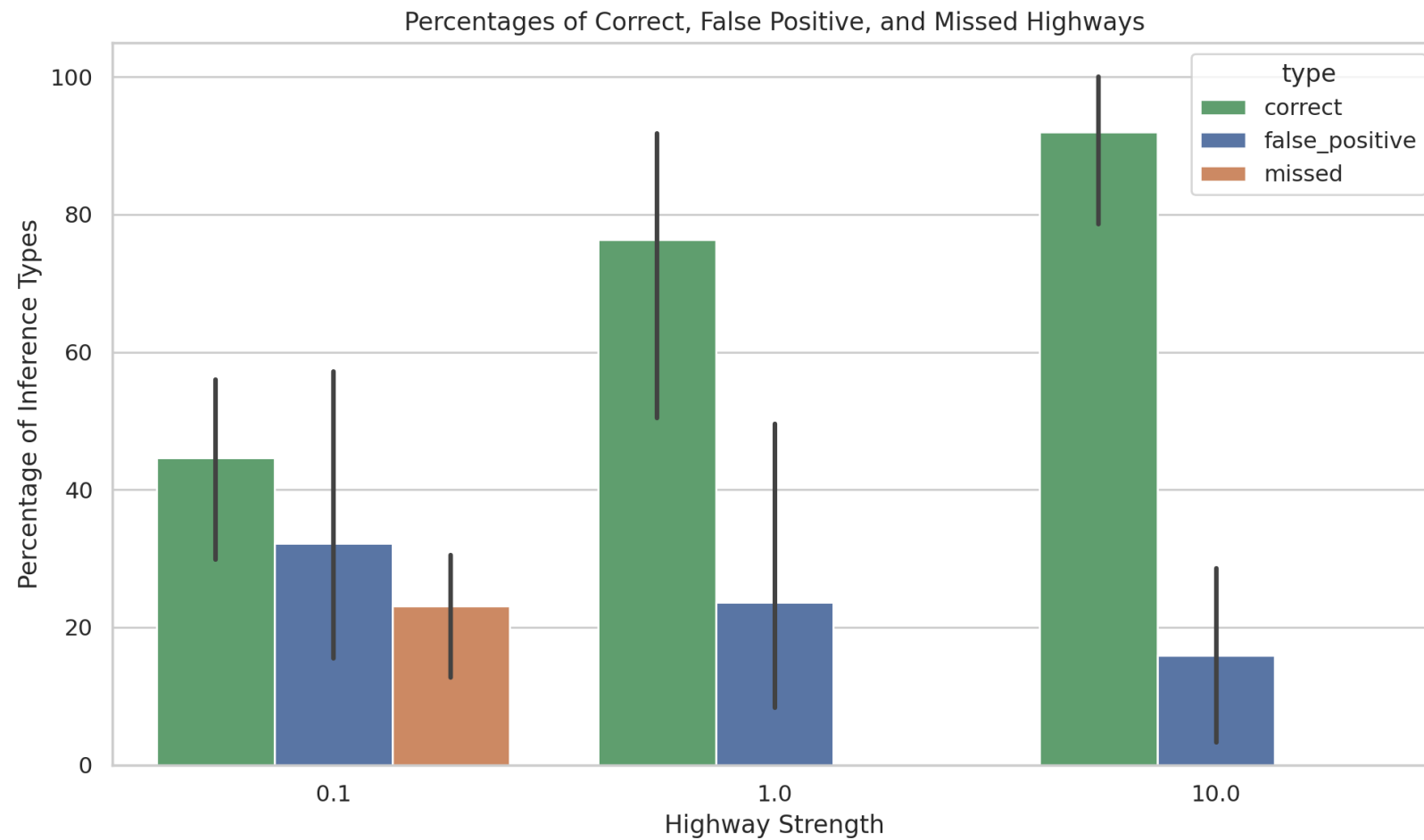


# Results (Transfer Heuristic)

$$\text{BIC: } \alpha \times \Delta LL > \log(N)$$

Weak pruning (distance, LL improvement, BIC  $\alpha = 2$ )

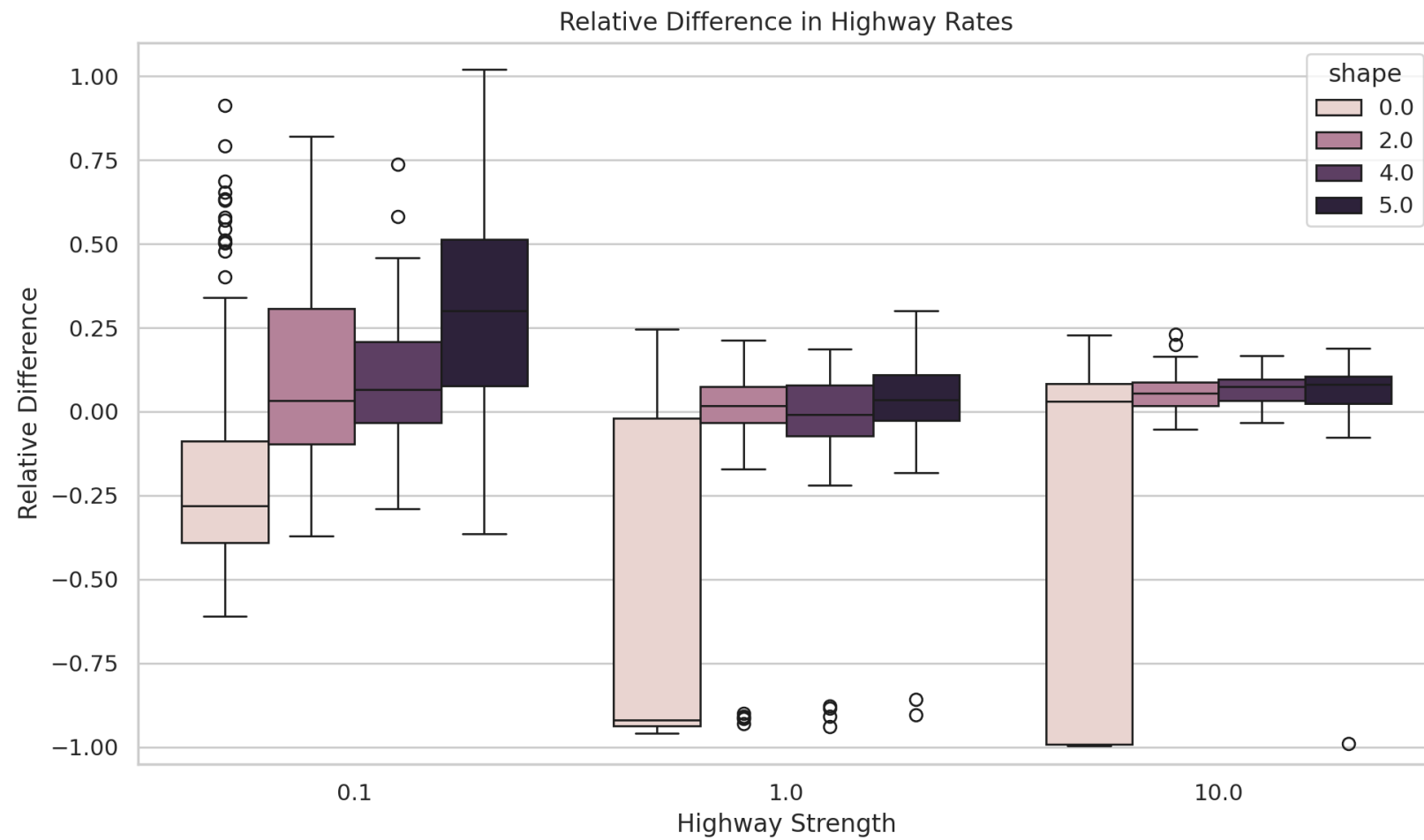
Strong pruning (distance, **potential highway rate**, LL improvement, BIC  $\alpha = 1$ )



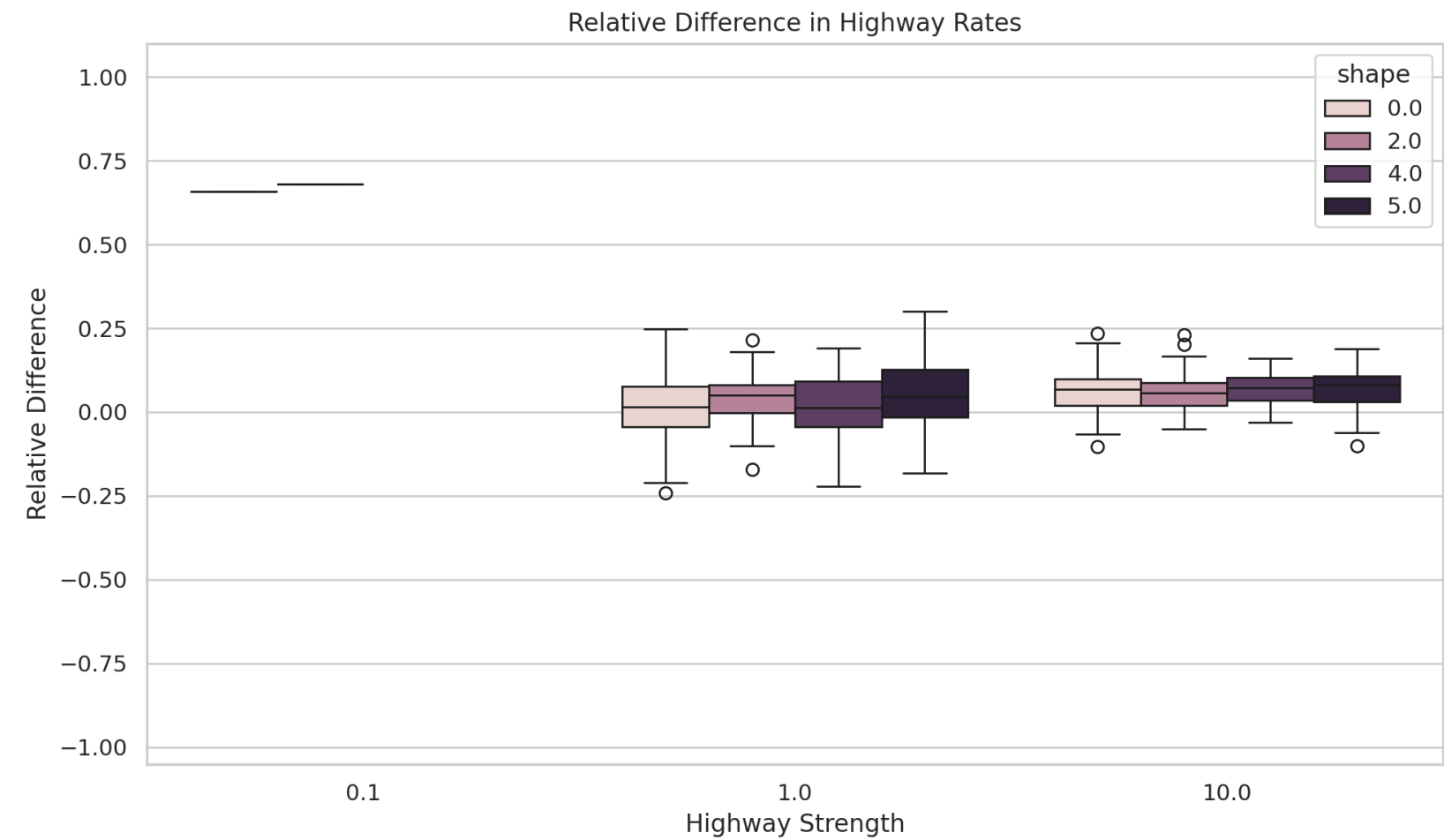
Depending on desired **sensitivity** to highway signals and/or runtime constraints heuristic may be adjusted

# Results (Transfer Heuristic cont.)

Weak pruning (distance, LL improvement, BIC  $\alpha = 2$ )



Strong pruning (distance, **potential highway rate**, LL improvement, BIC  $\alpha = 1$ )



# Limitations

- Missing metric of **realism** of simulated gene trees and their distributions
- Many other possibilities to model highways
- Empirical data is a lot more messy
- Empirical tests show many highways with varying strength, **how many should we capture?**
- Big **additional runtime** for the parameter optimization
- Some **numerical issues** remain

## Future Work

- (More) empirical tests
- Different ways of simulating gene trees and highways
- Different **penalties/optimization functions** to avoid overfitting
- Joint optimization of DTL + highway parameters
- **Faster inference** of highway rates that are "good enough"
- Investigating the impact of **ghost lineages** on highways

## Appendix: Likelihood Math

$$\Pi_{e,\gamma} = p_e^S \sum_{\gamma',\gamma''|\gamma} p(\gamma',\gamma''|\gamma)(\Pi_{f,\gamma'}\Pi_{g,\gamma''} + \Pi_{f,\gamma''}\Pi_{g,\gamma'}) \quad (S)$$

$$+ p_e^S (\Pi_{f,\gamma}E_g + \Pi_{g,\gamma}E_f) \quad (SL)$$

$$+ p_e^D \sum_{\gamma',\gamma''|\gamma} p(\gamma',\gamma''|\gamma)\Pi_{e,\gamma'}\Pi_{e,\gamma''} \quad (D)$$

$$+ 2p_e^D \Pi_{e,\gamma}E_e \quad (DL)$$

$$+ p_e^T \sum_{\gamma',\gamma''|\gamma} p(\gamma',\gamma''|\gamma)(\Pi_{e,\gamma'}\bar{\Pi}_{e,\gamma''} + \Pi_{e,\gamma''}\bar{\Pi}_{e,\gamma'}) \quad (T)$$

$$+ p_e^T (\Pi_{e,\gamma}\bar{E}_e + \bar{\Pi}_{e,\gamma}E_e) \quad (TL)$$

$$+ \sum_{d \in H_e} p_{e,d}^H \sum_{\gamma',\gamma''|\gamma} p(\gamma',\gamma''|\gamma)(\Pi_{e,\gamma'}\Pi_{d,\gamma''} + \Pi_{e,\gamma''}\Pi_{d,\gamma'}) \quad (H)$$

$$+ \sum_{d \in H_e} p_{e,d}^H (\Pi_{e,\gamma}E_d + \Pi_{d,\gamma}E_e) \quad (HL)$$

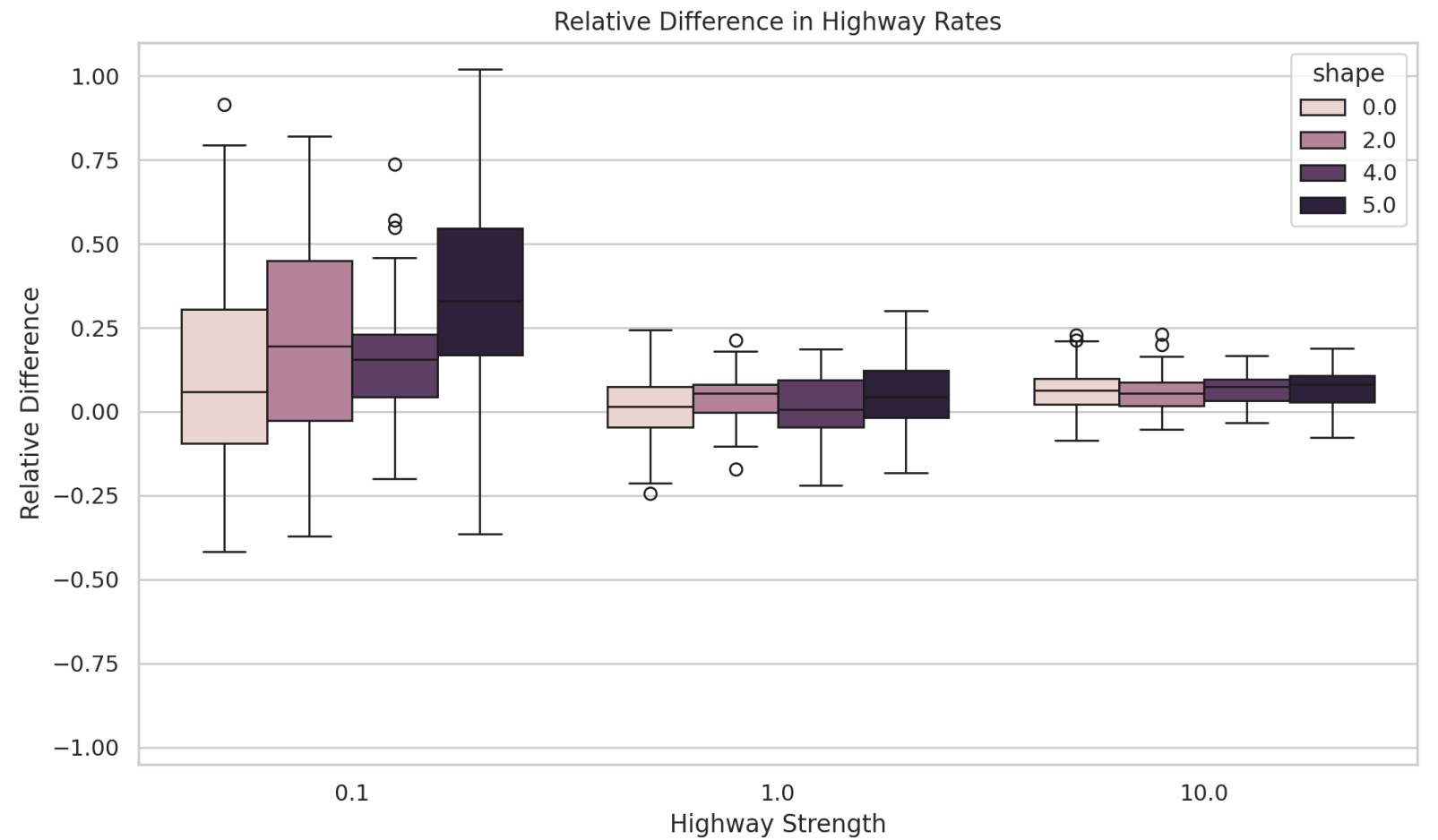
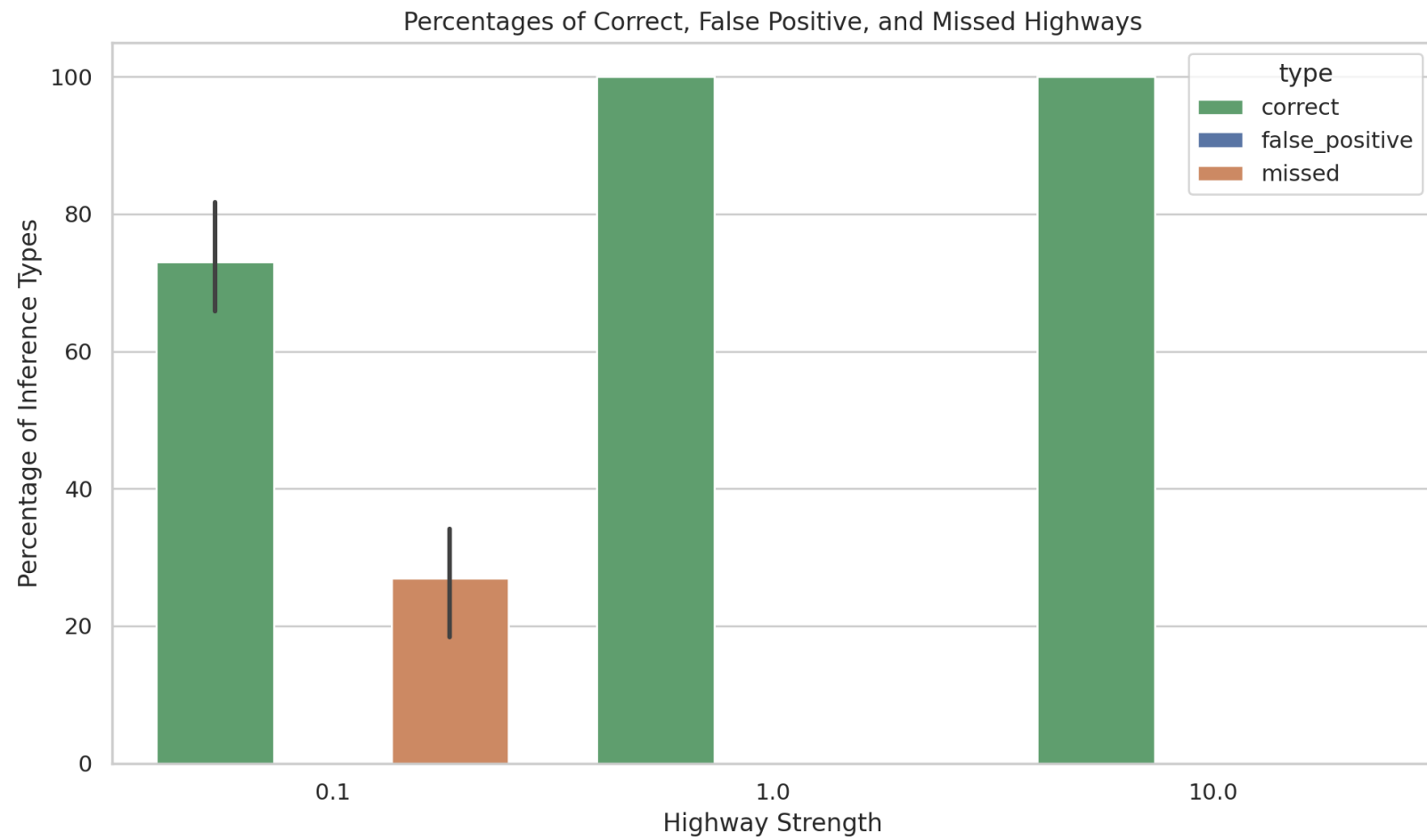
## Appendix: Likelihood Math (cont.)

$$E_e = p_e^L + p_e^S (E_f E_g) + p_e^D (E_e^2) + p_e^T (E_e \bar{E}_e) + \sum_{d \in H_e} p_{e,d}^H (E_e E_d)$$

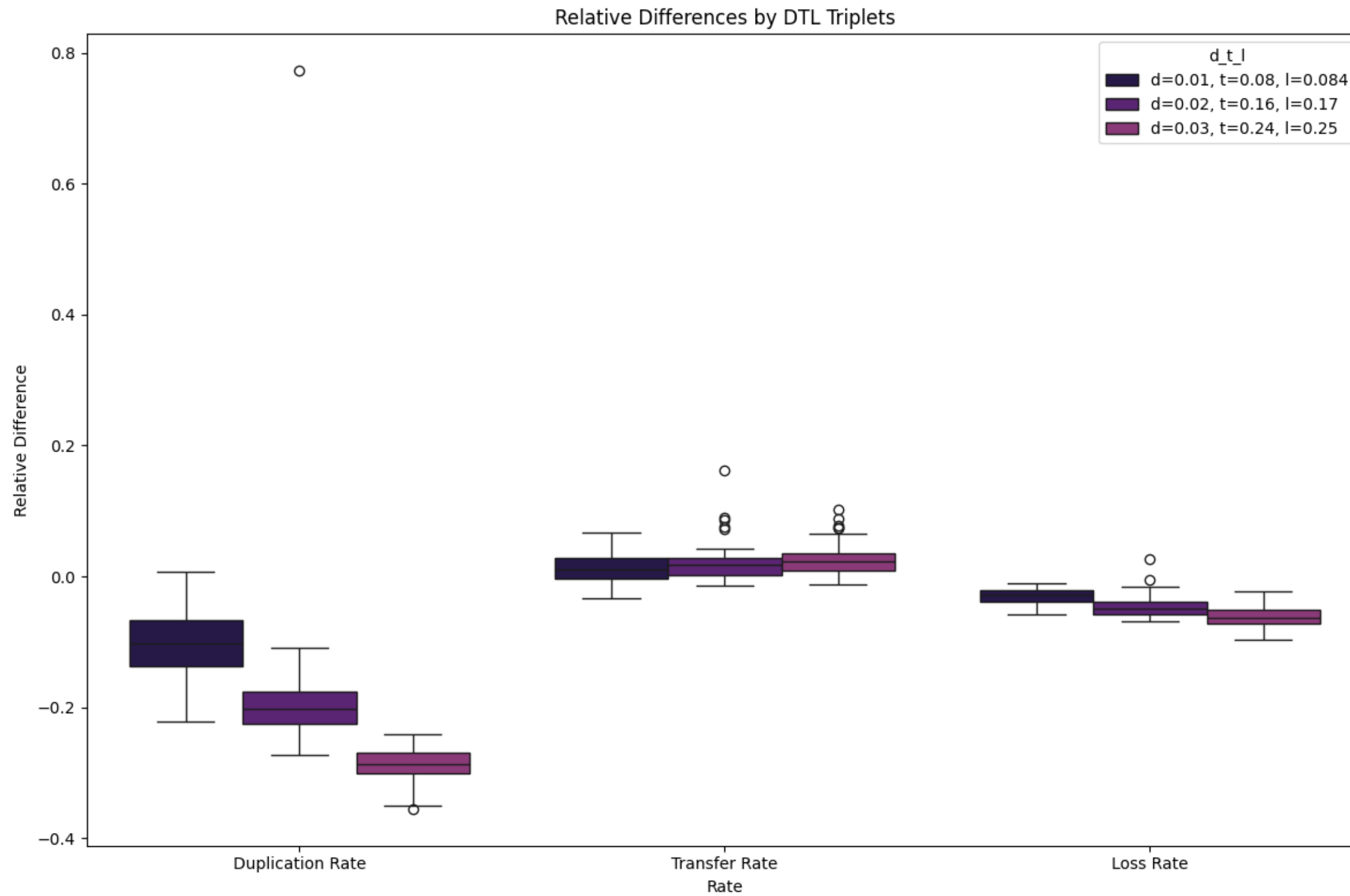
For a **gene alignment**  $A$  and a **species tree**  $S$ :

$$P(A|S) = \frac{\sum_e p_e^O \Pi_{e,\Gamma}}{\sum_e p_e^O (1 - E_e)}$$

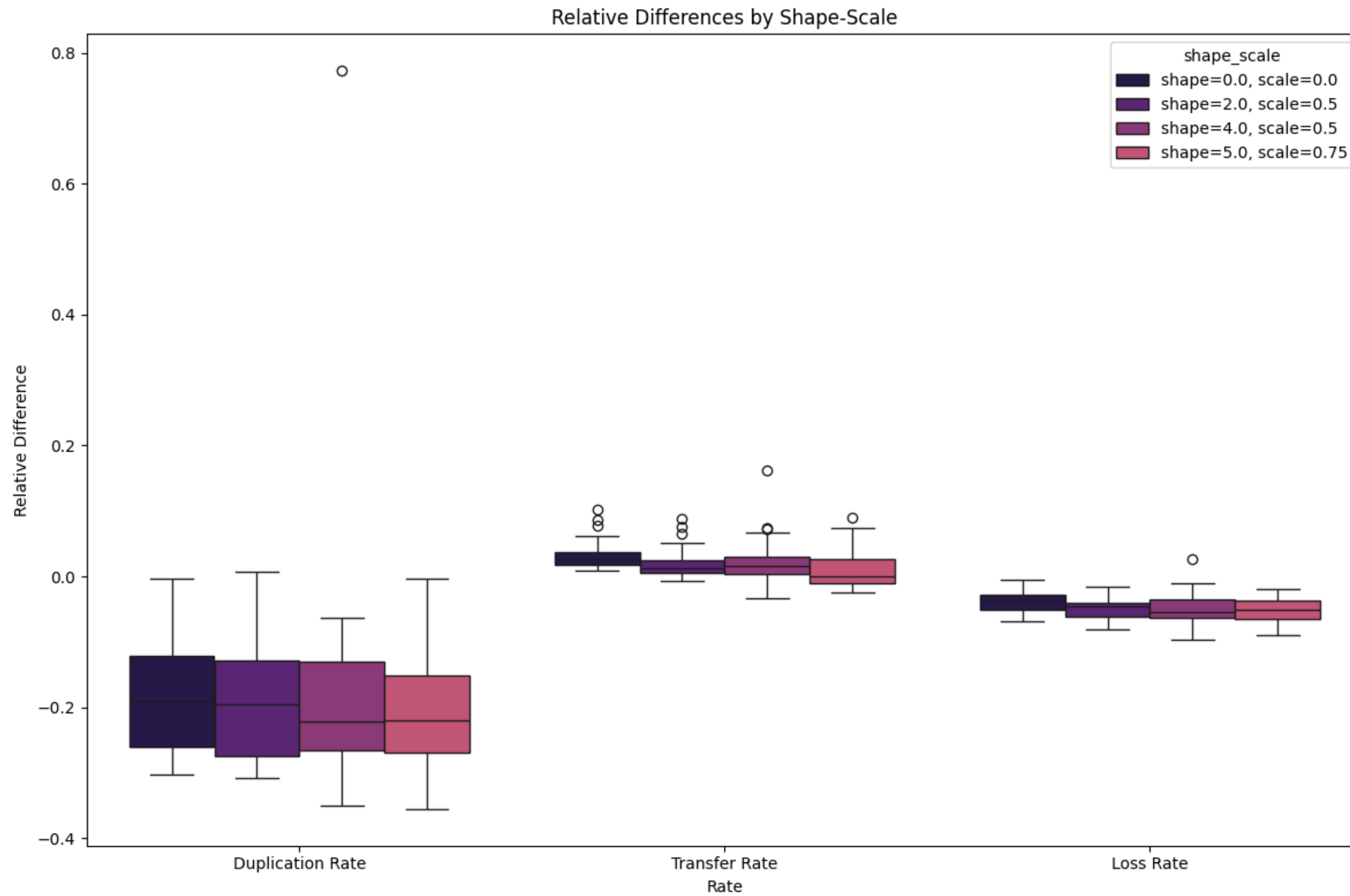
# Appendix: Results (Highway Validation)



# Appendix: Results (DTL)

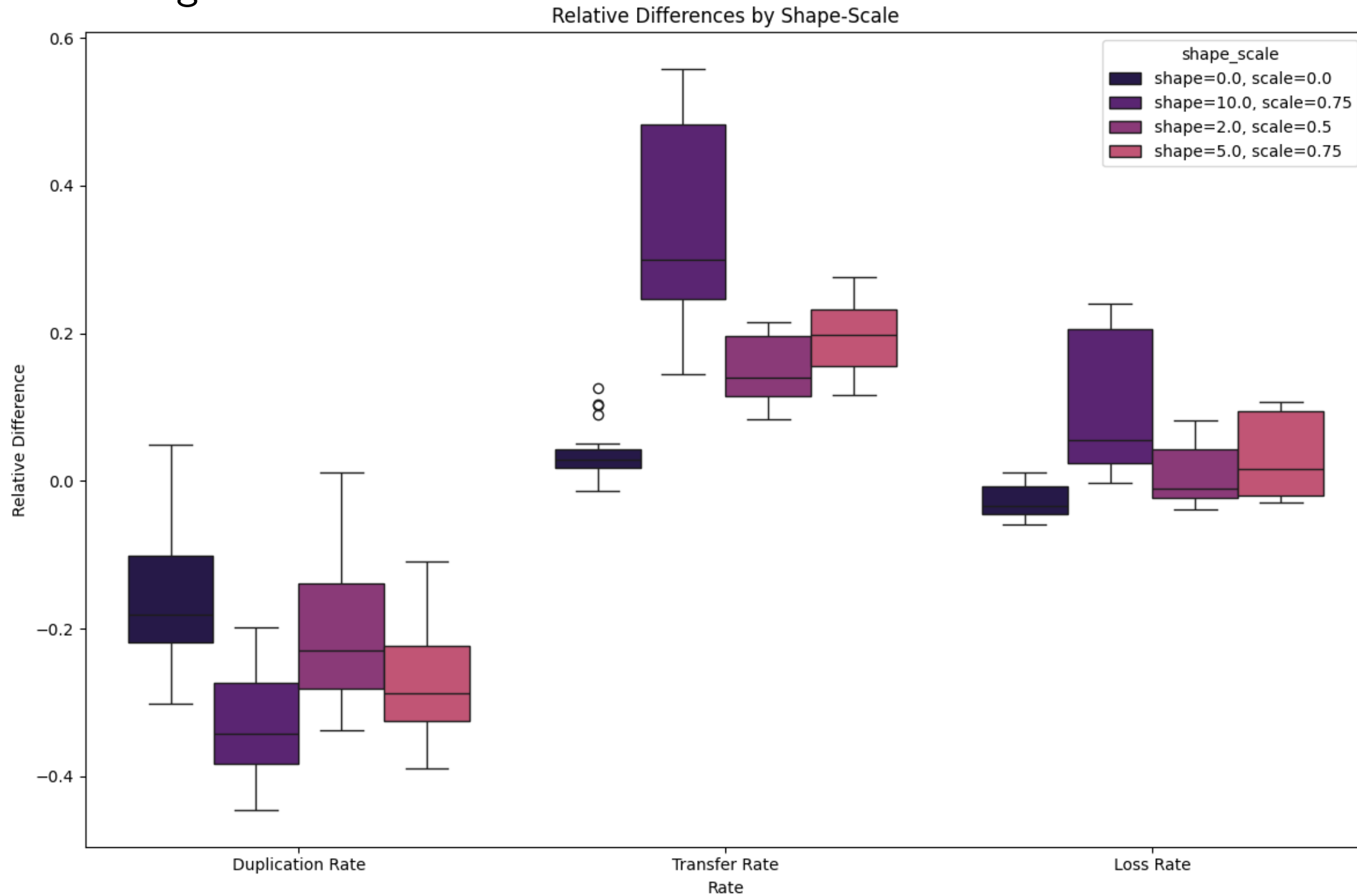


# Appendix: Results (DTL cont.)



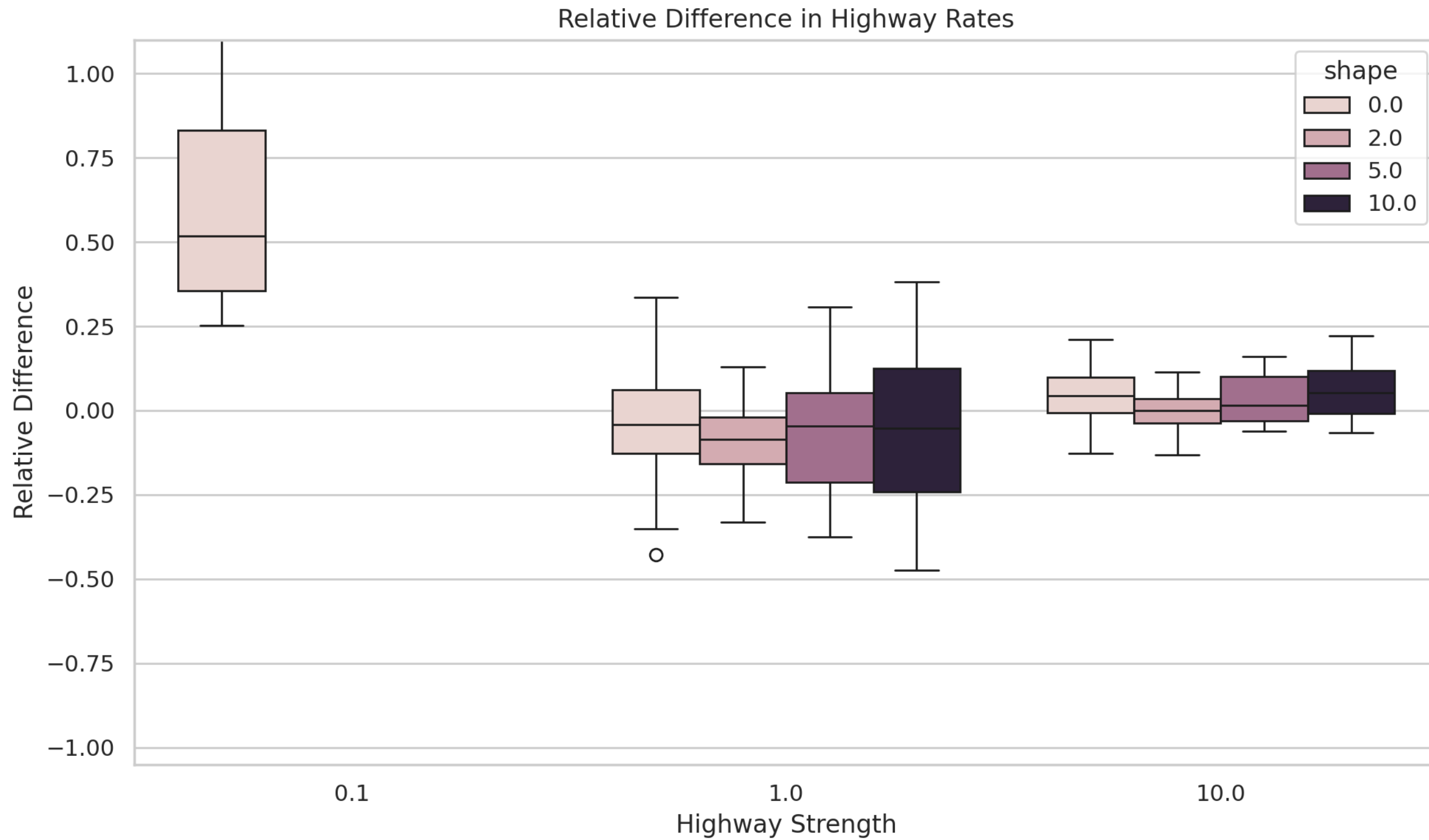
# Appendix: Pushing Error and Noise

Additional NNI moves on ground truth tree

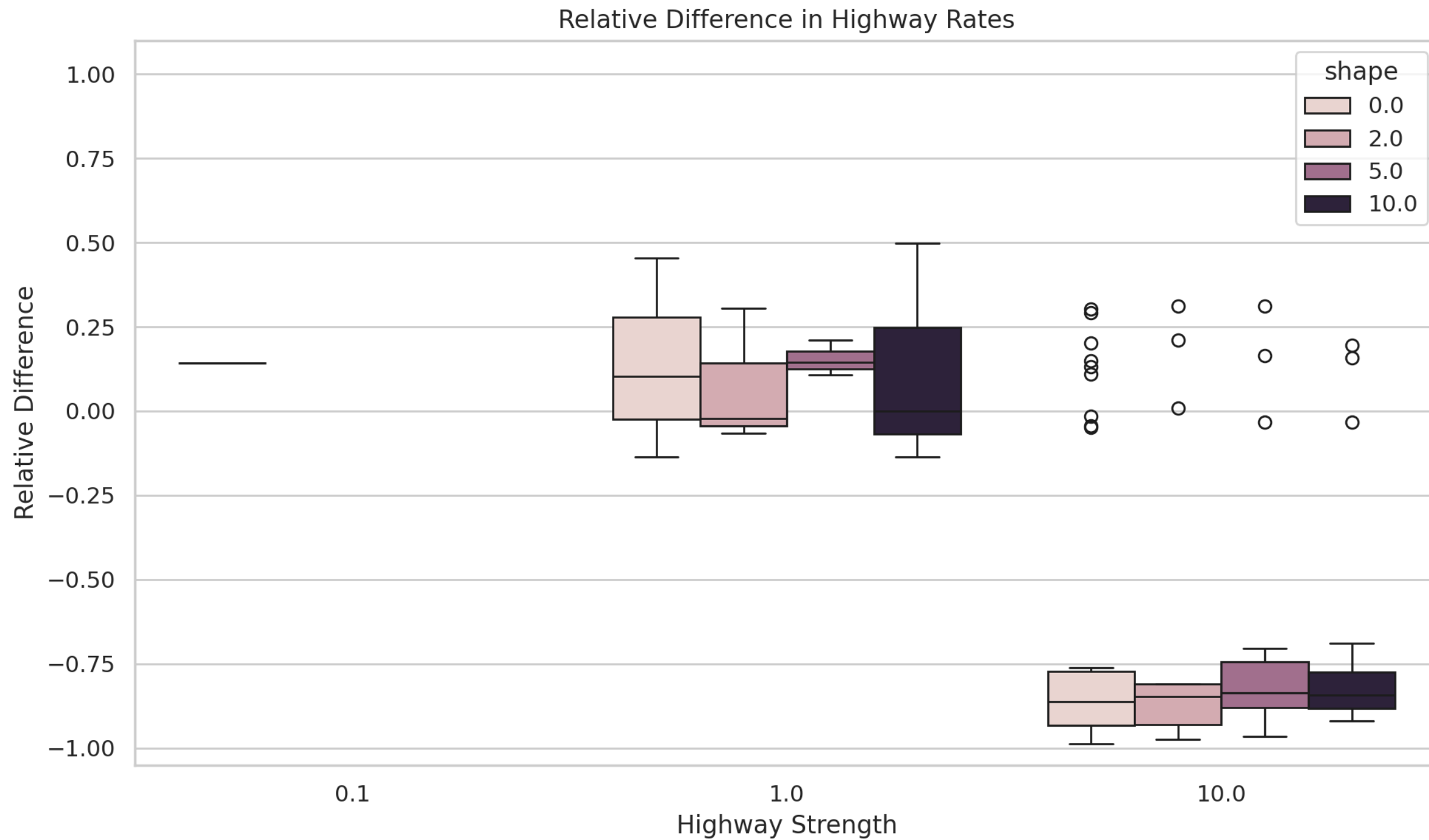


D, T, L = 0.1, 0.3, 0.4

# Appendix: Pushing Error and Noise



# Appendix: High Rates and Error



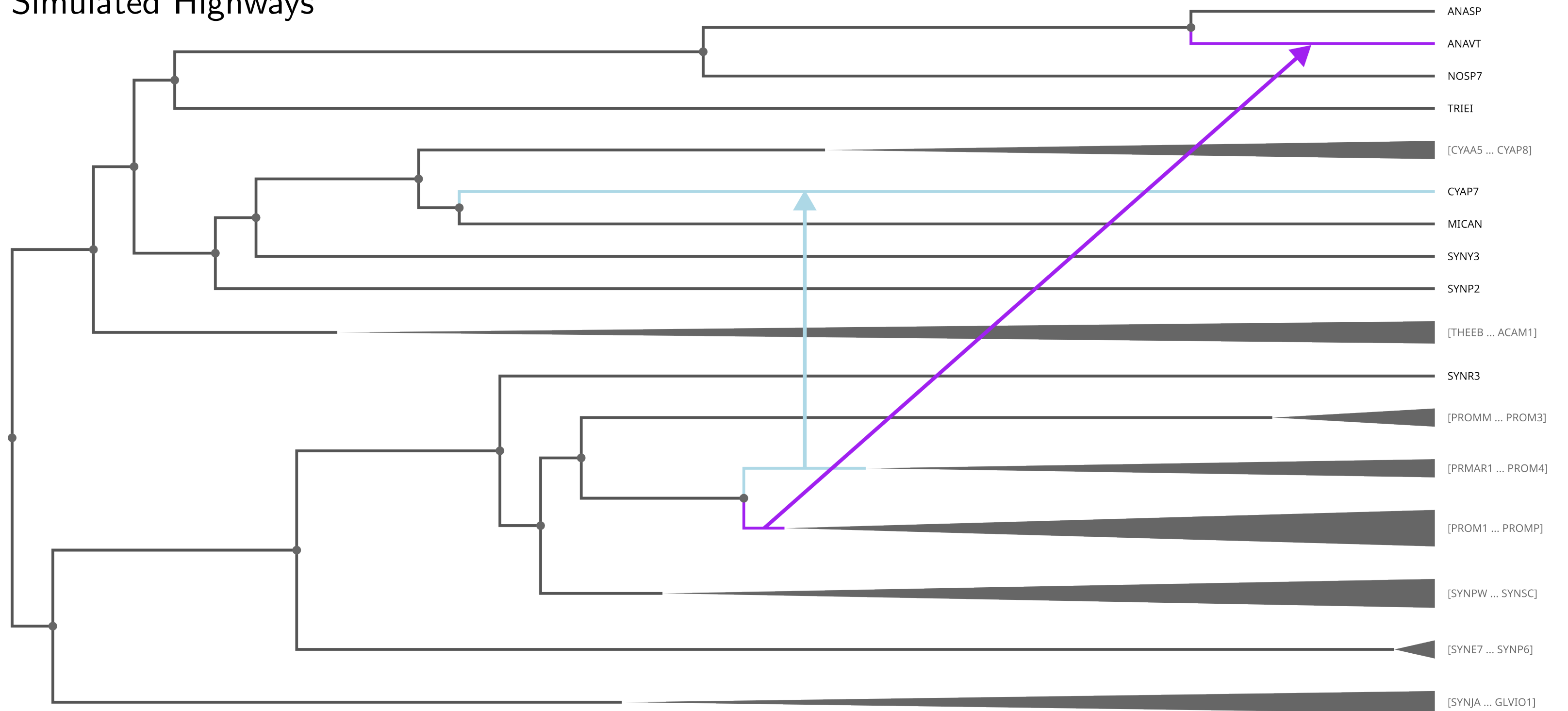
## Appendix: Results (Complex Simulations)

Change in DTL inference error: ( $\Delta D = 0.145$ ,  $\Delta T = 0.157$ ,  $\Delta L = -0.016$ )



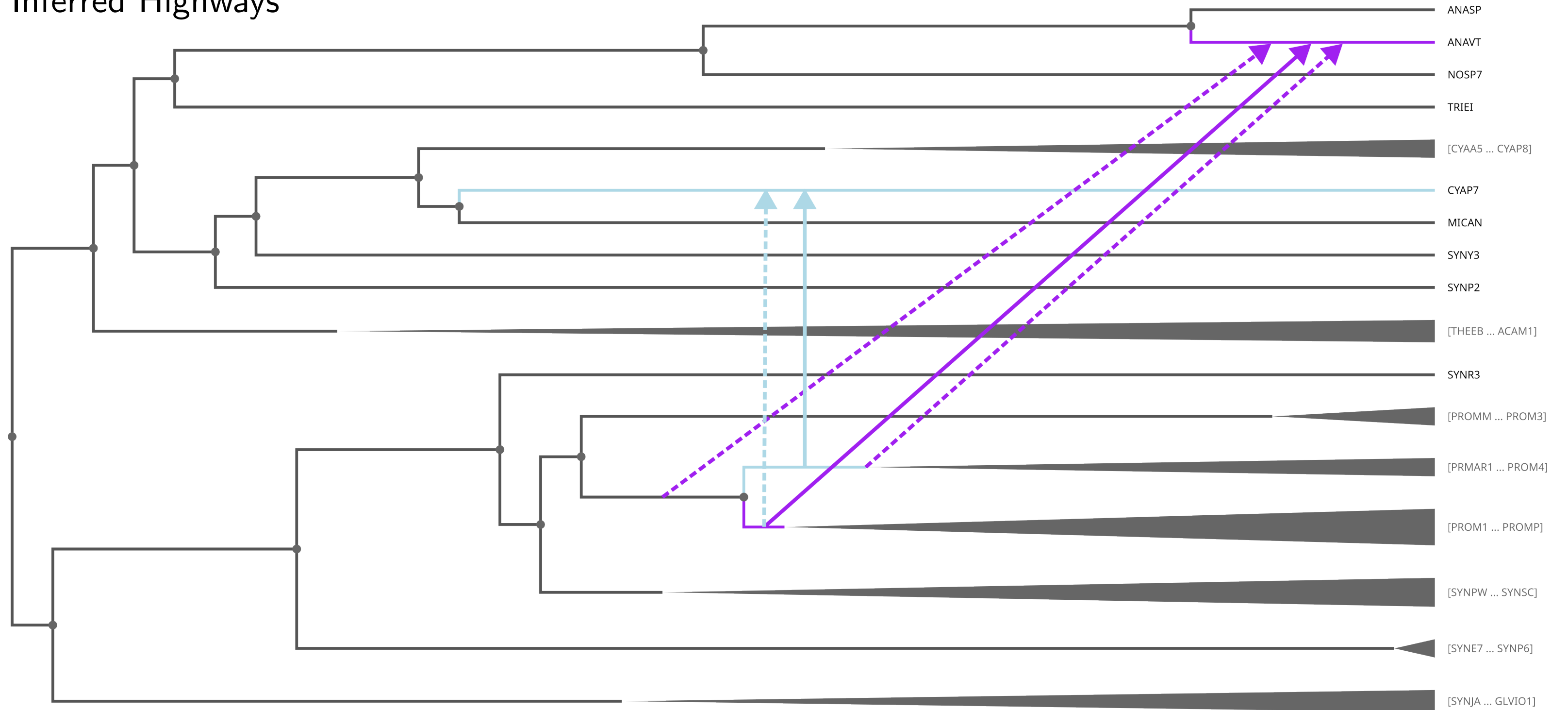
# Appendix: Results (Complex Simulations)

## Simulated Highways



# Appendix: Results (Complex Simulations)

## Inferred Highways



# Appendix: Results (Complex Simulations)

## Inferred Highways

