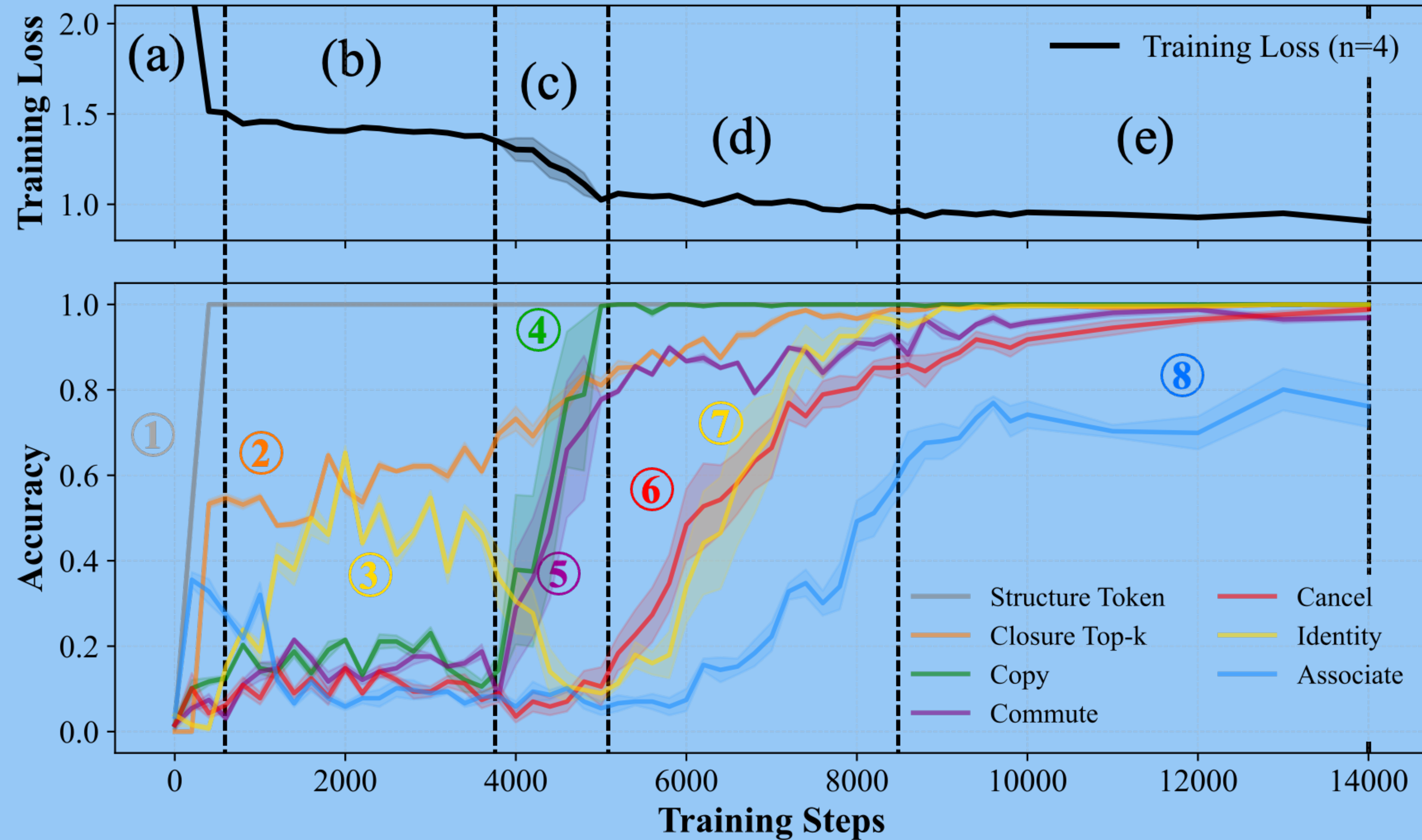


# Transformers infer structure from variables



## In-Context Algebra

Eric Todd, Jannik Brinkmann, Rohit Gandikota, David Bau

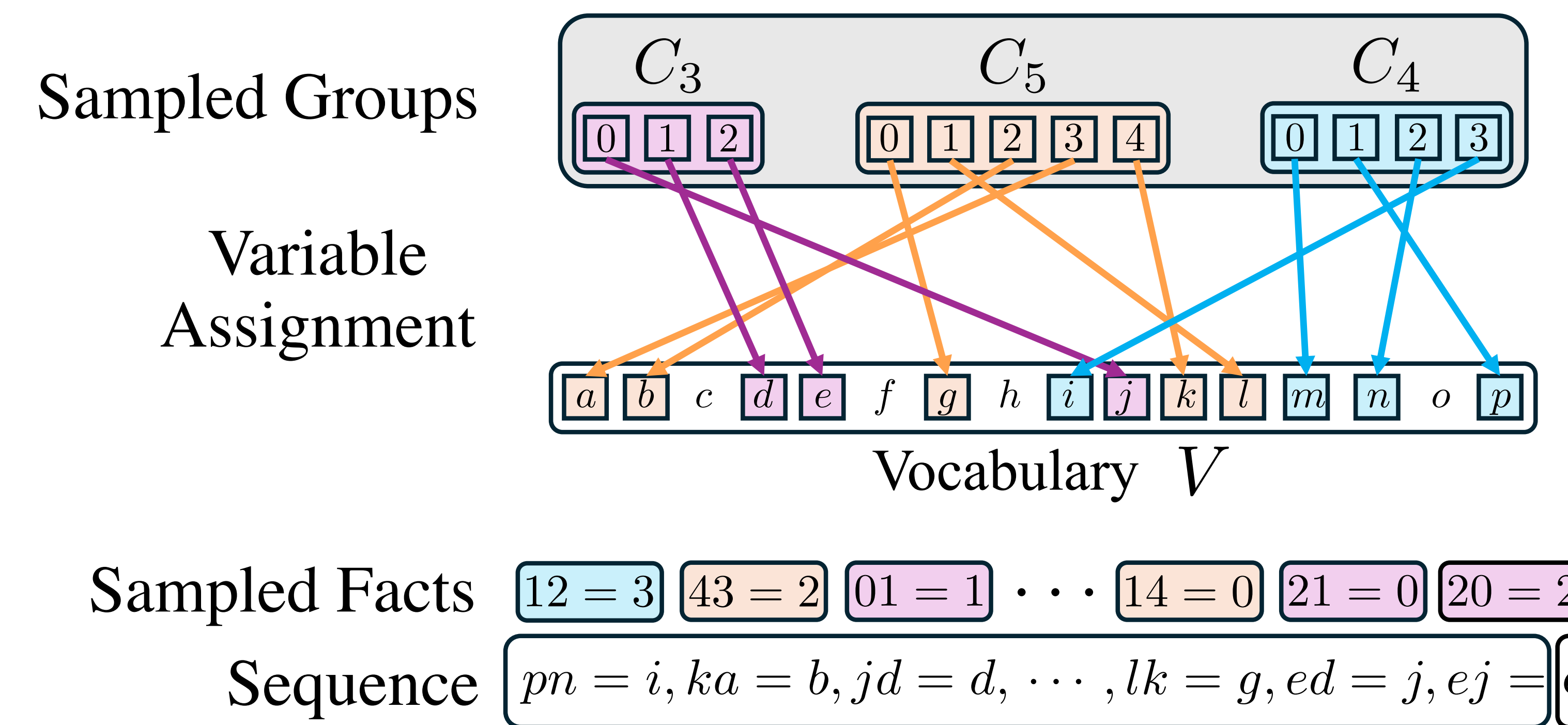


Website: [algebra.baulab.info](http://algebra.baulab.info)

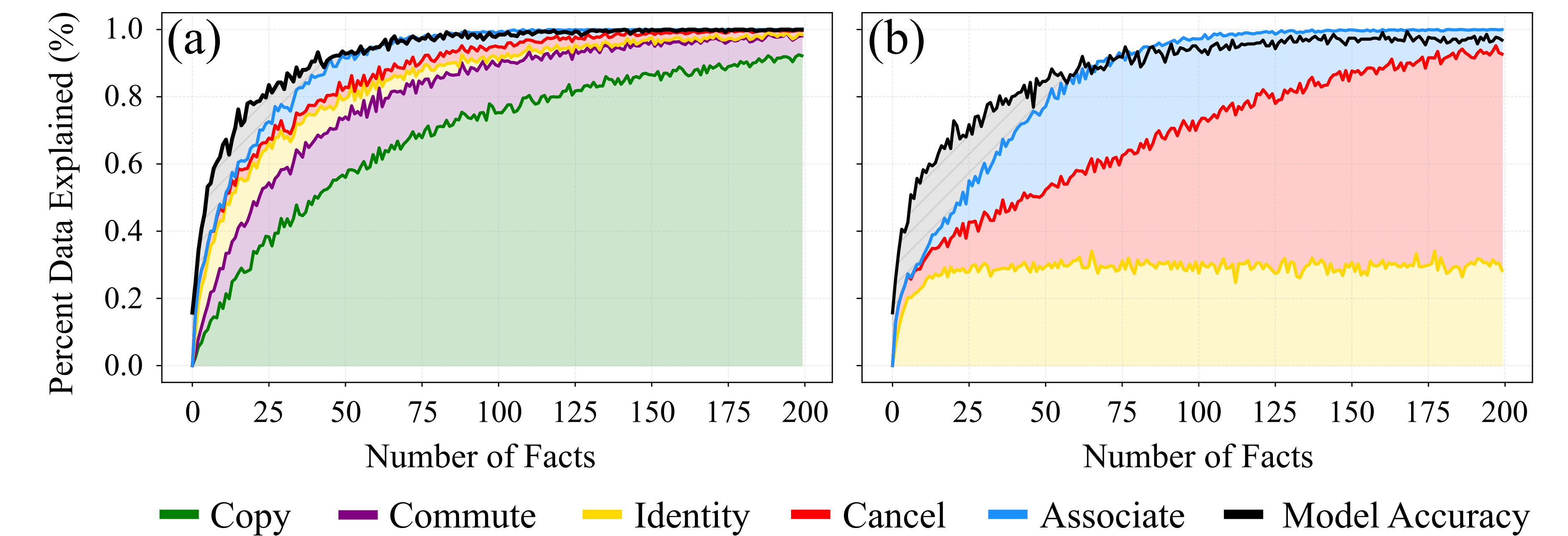


$12 = 3$   $43 = 2$   $01 = 1$   $\dots$   $14 = 0$   $21 = 0$   $20 = 2$   
 $pn = i, ka = b, jd = d, \dots, lk = g, ed = j, ej = e$   
 $30 = 3$   $41 = 5$   $60 = 6$   $\dots$   $24 = 0$   $21 = 7$   $13 = 4$   
 $bd = b, ai = c, gf = g, \dots, eh = d, li = k, pb = h$   
 $13 = 4$   $20 = 2$   $34 = 0$   $\dots$   $65 = 4$   $23 = 5$   $66 = 5$   
 $fp = k, ae = a, pk = e, \dots, lc = k, ap = c, ll = h$

(1) We design a **contextual reasoning** task where **interactions** define tokens



(2) We **isolate** parallel reasoning strategies using **targeted data distributions**



(3) Transformers develop discrete symbolic strategies to reason about variables

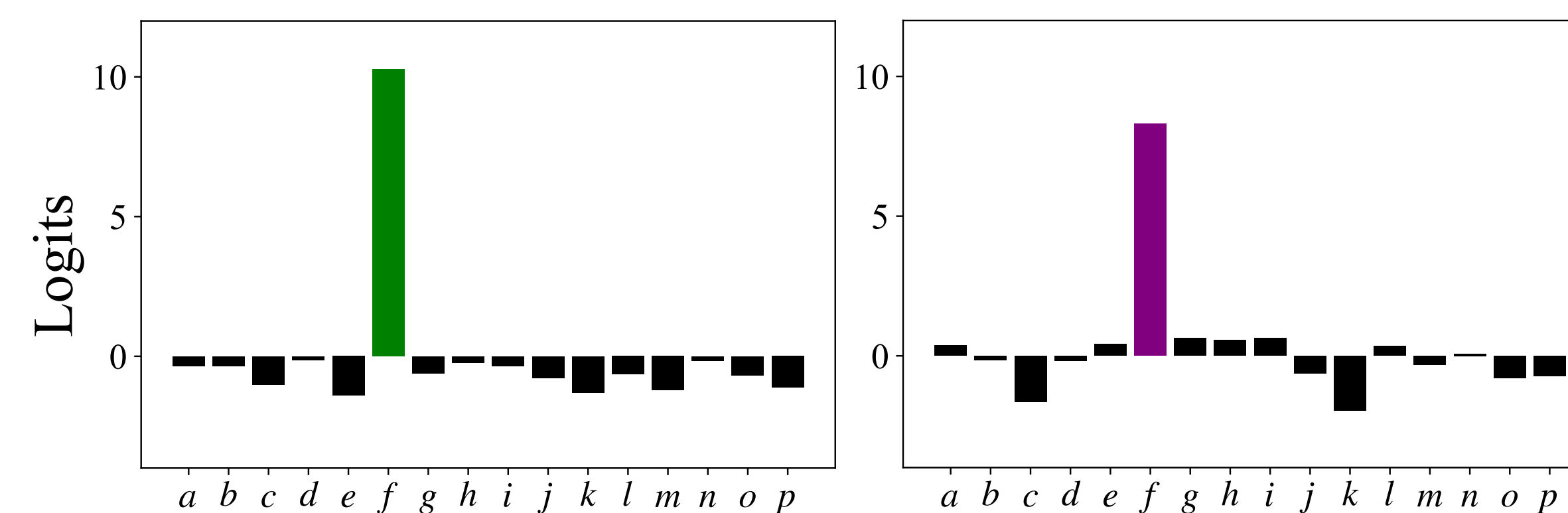
## Commutative Copying

$, go = j, kg = e, fo = e, jb = c, gj = f, eo = o, jo = k$   
 $, ck = f, og = j, gg = b, oc = g, ob = f, bj = c, kc = f$   
 $, eo = o, ke = k, bf = k, ge = g, jk = o, bo = f, jk = o$   
 $, cb = e, bc = e, kc =$

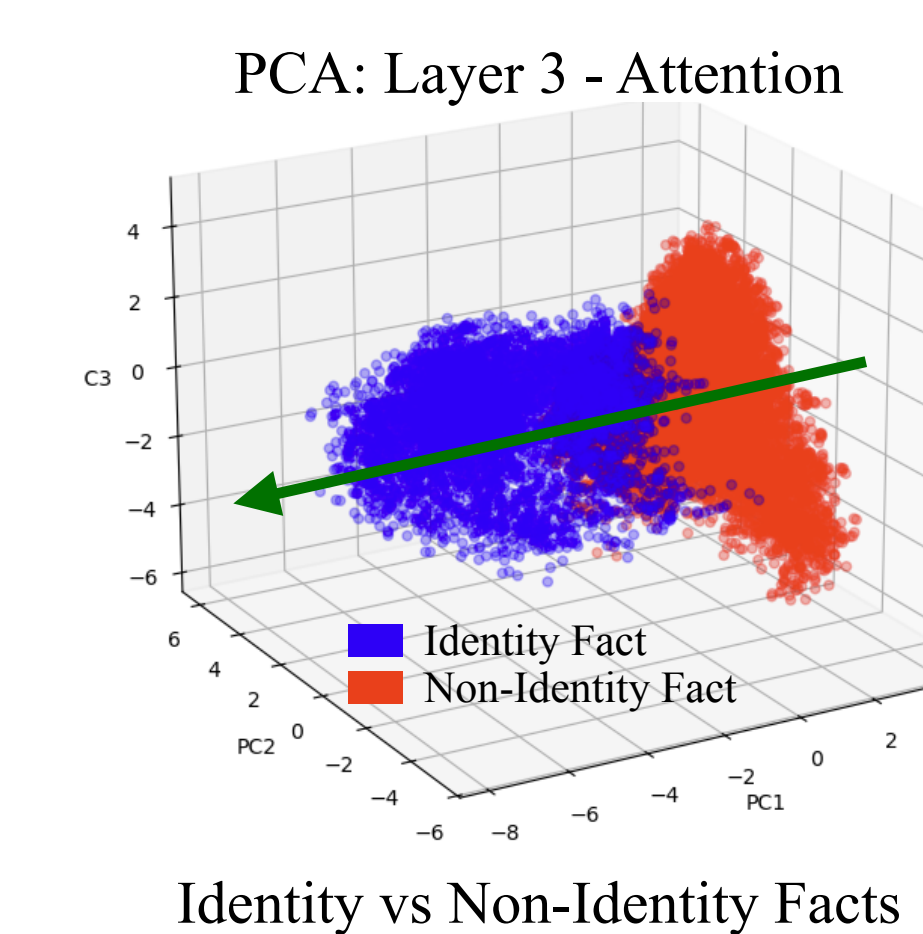
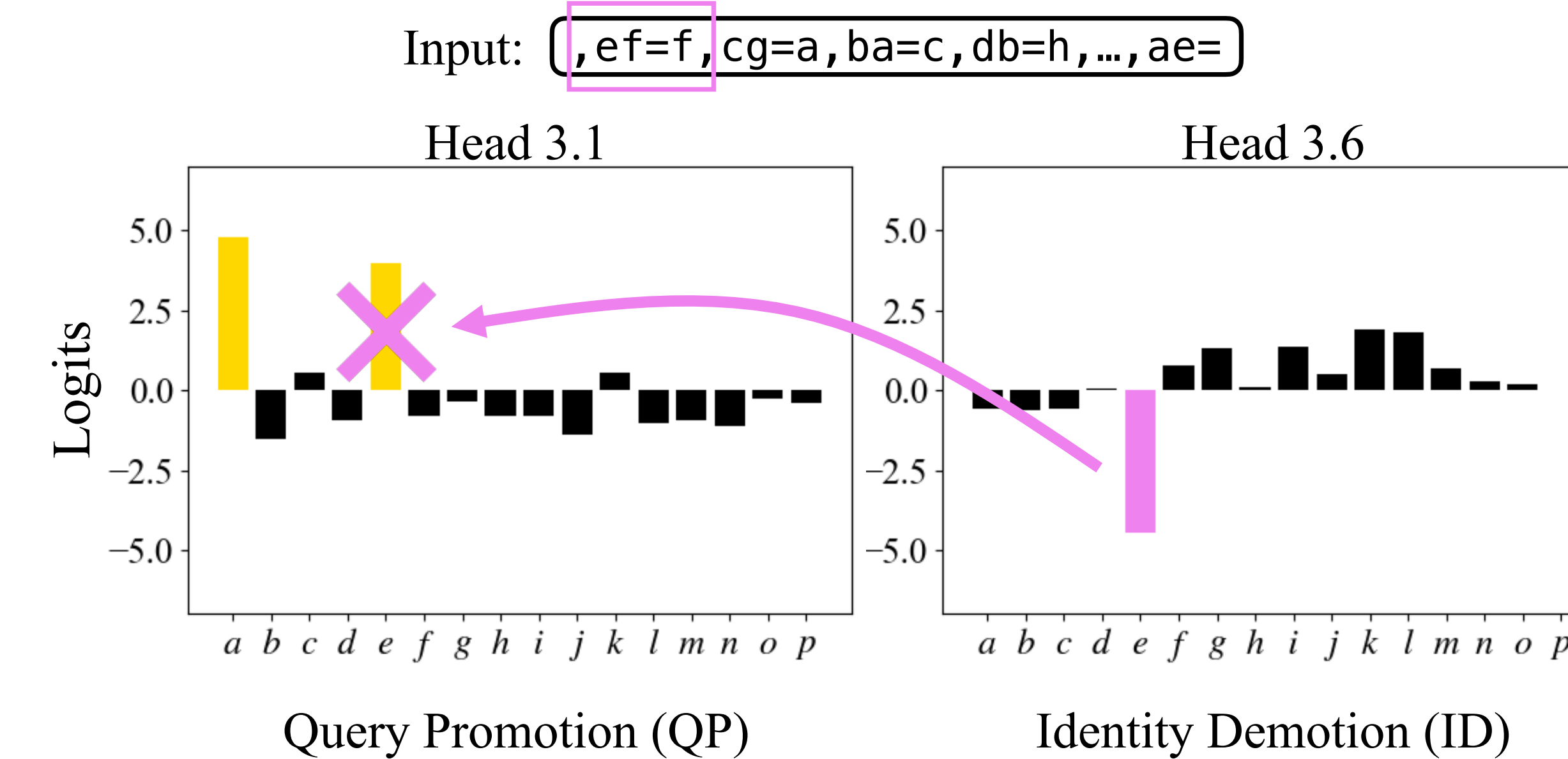
Verbatim Copying

$, go = j, kg = e, fo = e, jb = c, gj = f, eo = o, jo = k$   
 $, ck = f, og = j, gg = b, oc = g, ob = f, bj = c, oj = k$   
 $, eo = o, ke = k, bf = k, ge = g, jk = o, bo = f, jk = o$   
 $, cb = e, bc = e, kc =$

Commutative Copying



## Identity Recognition



- (i) Clean:  $, ad = f, cg = e, fa = c, db = d, \dots, df = e$
- (ii) PCA Steering (QP):  $, ad = f, cg = e, fa = c, db = d, \dots, df = d$
- (iii) QP + Demote d:  $, dh = h, ad = f, cg = e, fa = c, db = d, \dots, df = f$
- (iv) QP + Demote f:  $, fh = h, ad = f, cg = e, fa = c, db = d, \dots, df = d$

## Cancellation

