



# **Analysis of Transformer Behaviour in Reinforcement Learning**

**Miguel Silva Amaral**

Thesis to obtain the Master of Science Degree in

**Electrical and Computer Engineering**

Supervisor: Prof. Arlindo Manuel Limede de Oliveira

## **Examination Committee**

Chairperson: Prof. Pedro Filipe Zeferino Aidos Tomás

Supervisor: Prof. Arlindo Manuel Limede de Oliveira

Members of the Committee: Prof. Alexandre José Malheiro Bernardino

Prof. Francisco António Chaves Saraiva de Melo

**December 2022**

# Acknowledgments

CS&gt;R... ^z zb zP- ^WzPCeqfKssbqzP- z- Yb..C@zPS ..bqWzb 4C-b\ CebssSYC>eqfHl, qS^@b a YfCSq  
..Pb eY<C@PS zq-sz S \ C-^@PCY@C\ C...SP PS f- Y- 4CS'e-z S' zPC \ -VSL bHzPS zPCsSi  
rCb^@R... ^z zb zP- ^W %H\ S%HqzPCsq YfC-^@s-eebqz bfCqzPC %G q>..SPb-z zPC\ R..b-Y@  
P- fC^CfCqP- @-^ beebqz-^S%zb •b-qfSP -s R@S^S' zPCYsz " fC %G qsi  
X- sz %zb \ %HfC^@>zP- ^W Hq- Yb..S'L \ Cz b Lq... - Y^LsSC %q- qfCsi R<- ^ z YbW4 <W- z zPS  
U-q^C%a..SPb-z HCS'Leq~@bH..P- z...C- <<b\ eYSPC@ RWb...zP- z %qHz-qC...SY4C4qLPzi Rz..b-Y@  
4C- @ssCqfSCzb %qHfC^@PS e zb ^bz \ C^zb^ %q^- \ Gs, ^z ^S>3G zqf di>3G zqf pi>3Cq^- q@>  
; -z qf^->? SIlb>OC^qf ~C-^@yb\ -si



R@C<YqC zP-z zPS @b<~\ C^z S - ^ bqLS-Y..bqWbH\ %b..^ --zPbqPSe - ^@ zP-z S HY Y - YzPC  
qd ~Sq\ C^zs bHzPC; b@CbH; b^@~<z - ^@Kbb@dq <zCs bHzPC} ^SfCpS@ @C@XS4b-i



# Abstract

? C<S^ yq ^shh\ Cqf? yg9: S'zpb@<Gs sD ~C^<C\ b@CS^L - s - \ CzPb@HhQ- <PSCfS^L @GSp@b-z\b\ Gs S' qS^HhQ\C C^z YG q^S^L fp Xg eqp4Y\ s> ..PSC - fbS^L <b\ \ b^ Ss-Gs zP-z <- ~sC @SfCqL^<C S^ zPS sCzS^L 9:i rz-zGs>- <zb^s - ^@Hz-qC qz-q^s - qC-s\ eY@Hh\ - @-z-sCz S' bq@Cq zb <b^@Sb^ Pb...zPC \ b@CY..SYeY^ - seG^S <zq UC-zbq^..PSC \ -sW^L Hz-qC e- qz bHzPC sD ~C^<C zb \ - WC G <P zS GszCe -zzC^@ b^%zb Ss e- sz S' - ^ - ~zbqLqGssSfC HsPSb^i rb\ C bHzPC ? y s 4bzz^C G^W - qC S@C^zS C@ S' zPS ...bqW - ^@- q^PSG-z-q Y<P- ^LGs - qC S\ eY\ C^zC@S' bq@Cq zb sz- 4SS C Pb...zPS - YbqzP\ eGqHh\ s 4Cz..CC^ G <P CeSb@C bHzPC C^fSp^ \ C^zi 3G<- ~sC ? y S S^ S@C 4%zPC @SfCpS^% ^@ l ~- Y^%bHSs @ z-sCz>4bzP - ^ b' QbS^% ^@- ^ b^QbS^%zq ^shh\ Cq- YbqzP\ ..CqC@CfCbeC@- s eqpbHhQ^<Cez bH Pb...CfeYbq zb^ ..b-Y@- eeY%S^ sD ~C^<C\ b@CS^L - YbqzP\ si GS^ - Y% sPSH Hh\ ~s^L zb%eqp4Y\ s zb..bqW^L...SP qC Yq.bq@" ^- ^<S Y\ - qWz ...: s S\ eY\ C^zC@zb sPb...zP-z pX<- ^ - Yb 4C~sC@zb sbYfC eqp4Y\ s zP-z b<<-q S' zPC qG Y..bq@ BfeGp\ C^zs @C b^szq zC zP-z zPCqC S\$ qb\ HhQ S\ eqp4Y\ C^z S' <-qC^z sD ~C^<C\ b@CS^L sbY-zb^s>- ^@ zP-z zq ^shh\ CqP P- fC- sL^S <- ^z - \ b-^z bH-^z eeC@ ebzC^zS Y..PC^ zPC% qC ~sC@S' @CeC^@C^zY%HhQ CfeYbq zb^i

# Keywords

? Cce pCS^HhQ\C C^z XG q^S^L yq ^shh\ Cqst rD ~C^<C [ b@CS^L Lt GS^ - ^<S Y [ - qWzst [ - qWf ? C<S^ dqb<GssGst



# Resumo

? C<Sb^ yq ^shbq\ Cqf? yg9: S'zpb@< \ b@CY éxb @CsD ~ ^<S s<b\ b~\ - \ Cz@bYLS e-q - Y- ^é qbs  
qS~ Y- @bs @CsU @bs C\ eq4Y\ -s @C- eq^@S- LC\ ebq qHhqb>C^l ~ ^zb CfSz- eq4Y\ - s<b\ ~^s @Csz-  
-eG >l ~Cz&S<- \ C^zC< ~s- \ @SfCqL ^<S ^- -eq^@S- LC\ 9:i Bsz @bs>- é Cs Cq<b\ eC^s-s Hz-q s sxb  
- \ bszq @bs @C~\ @z sCz @CHbq\ - - <b^@S-Sb^-q- Hbq\ - <b\ b~\ \ b@Cb eY^CS ~\ - <Cqz- zq UZ fS >  
\ -s<- q ^@b sG<é Cs Hz-q s @CHbq\ - - l~C< @ S'zCqf- Yb zC\ ebq Y- eC^s <b^sS@CqC b sC- e-ss- @b @C  
Hbq\ - - ~zbQ(LqCssSf- i | CszC zq 4- Ypb sxb S@C^zS < @bs - Y~^s eq4Y\ - s S'Gq^zCs - b? y C\ ~@ ^é-s  
- d ~SzC-q S Hbq\ \ S eY\ C^z @s @CHbq\ - - Csz 4SS- q- Hbq\ - @C<b\ b CszC\ b@Cb S'zCq LC- < @  
S'sz^a^<S @b eq4Y\ - i y- ^zb ~\ - Ybqfz\ b b^@bY%<b\ b b' @bY%<Hbq\ \ @CsC^fbYfS@bs @CHbq\ - -  
eqf-ql ~CCf eYbq éxb eb@CsCq- eY- @ - \ b@CY éxb @CsD ~ ^<S s>qsbYfC^@b- Y\ Sz- éxb @Cl ~ Y@ @C  
C @SfCqS@ @Cl ~C b ~sb @C~\ @z sCz <~s- i GS^ Y C^zC> HbS Cqz~ @ ~\ - zq ^sSéxb @C eq4Y\ - s  
CfC eY < zSfbs e-q \ Cq<- @bs " ^- ^<CSpS @C\ - ^Cqf - @C b^szq ql ~C- - eY- 4SS@ @C@CszC\ b@Cbs-  
eq4Y\ - s qG S ebss fCY- zq f s @C- Ybqfz\ bs @C- eq^@S- LC\ ebq qHhqb Cl ~Cb ~sb @Czq ^shbq\ Cqf  
~\ - e-ssb sL^S < zSfb ^- qsbY-éxb C' <S^zC @CszCs eq4Y\ - si

# Palavras Chave

, eq^@S- LC\ ebq pCHhqb dqbH^@ t yq ^shbq\ Cqf [ b@CY éxb @C rD ~ ^<S t [ Cq<- @bs GS^ ^<CSpst  
[ - qWf? C<Sb^ dqb<CssCst









I ; b^<Ysb^s - ^@G-z~qC,, bqW	J_
Iic , <PSCf\ C^zs	J_
Ii  G-z~qC,, bqW	IE
<del>3S4YbLq</del> eP%o	Ic
, ? C<Ssb^ yq ^sHbq\ Cq f- Y@ zsb^ qCs~Ys	Iu
, ic ObeeCqB^fSp^\ C^zpGs~Ys	ID
, i  „ -Wq @B^fSp^\ C^zpGs~Ys	v



# List of Figures

ic B^<b@Cq-qPSzCz-qC	ii	v
i  ?Cb@Cq-qPSzCz-qC	ii	u
i{ 3-sCyq^sHbq\Cq-qPSzCz-qC	ii	_
iJ , SsB^yq^sHbq\Cq, qPSzCz-qC	ii	_
iI [-sV@, zC^zB^	ii	cE
iv yq^sHbq\Cq X[ C b%ai	ii	cc
iu ; b\ e-qsB^ 4Cz.CC^kQC q^S'L-^@?CCe kQC q^S'L	ii	cJ
iD Kyqf X-qPSzCz-qC	ii	cu
i_ ?CSB^yq^sHbq\Cq, qPSzCz-qC	ii	cD
icE; b3BpX, qPSzCz-qC	ii	c_
icc ; b^zqszsC Yss S^ ; b3BpX	ii	c_
{ic O-YQ PCCz-P B^/Sp^ C^z	ii	{
{i  ?b-4CpGeY%3~' Cq, qPSzCz-qC	ii	D
{i{ yq^sHbq\Cq, qPSzCz-qC Hbqk QC q^S'L	ii	_
{iJ yq^sHbq\Cq, qPSzCz-qC Hbdqf;S -YdbS%a ez\ S-zB^	ii	{
Jic ; b^zCz XC^LzP=[ G^ pC.: q@-^@rz^@ q@? CfS zB^ -<qss O-YQ PCCz-P? -z-sCzs	iiiiiiiiii	{u
Ji  y-qLz [ b@CS'L -<qss O-YQ PCCz-P? -z-sCzs	ii	{D
Ji{ a fCq' zS'L 4b~^@ q%S [ C@S-\ Q†eCq @z-sCzs	ii	{_
JiJ dBSB^ -YB\ 4C@S'L=[ G^ pC.: q@-^@rz^@ q@? CfS zB^ -<qss O-YQ PCCz-P? -z-sCzs	iiiiiiiiii	JE
JiI K-zS'L=[ G^ pC.: q@-^@rz^@ q@? CfS zB^ -<qss O-YQ PCCz-P? -z-sCzs	iiiiiiiiii	J
Jiv ; -qebCpCs-Ys Hbqyq^sHbq\Cqk QC q^S'L	ii	J{
Jiu [ S'Lf@pCs-Ys Hbqyq^sHbq\Cqk QC q^S'L	ii	JJ
JiD K-\ Cszbe yq @S'L pCs-Ys Hbqyq^sHbq\Cqk QC q^S'L	ii	JI
Ji_ ; -qebCpCs-Ys Hbqyq^sHbq\Cqdda	ii	Jv
JicE  S'Lf@pCs-Ys Hbqyq^sHbq\Cqdda	ii	Jv



# List of Tables

{ic yq ^shdq Cqk Cq q^S'L P% Cq- q \ Cz Cq i {E

{i| yq ^shdq Cqdda P% Cq- q \ Cz Cq i {J

Jic pCs~ Ys b^ zPC,, - YCq? ? b\ -S' i {v

Ji| pCs~ Ys b^ zPCObecq? b\ -S' i {v

Ji{ pCs~ Ys b^ zPCO- YC PCCz- P ? b\ -S' i {v

JiJ pCs~ Ys ; b\ e- qsb^ bHa ^S' CpCs^ Hq-C\ C^z XG q^S'L ..S P rD ~C^<C[ b@CS^L i i i i i JD





# 1

## Introduction

, ~zb\ -zb^ P-s \ CqL@ -s b^C bHzPC HszGz Lq. S^L S^@-szqCs S^ qG-C^z %G qp> s^<C Sz C^ 4Ys 4-sQ  
^GssGs zb S^<qG sCzPCqC <S^<%<-z <bszs - ^@S\ eqj/CzPCqj sGj/S-G R^zC^L^z - LC^zs ^CC@ zb S^zCq <z  
..SZP zPCq C^fSj^ \ C^z b^ - @S^%4 sS - ^@..SZP qS^HqC C^z YG q^S^L zPC%< ^ S\ eqj/C <b^zS~b-s^%  
- Yb..S^L Hq- sPS^ S^ CjCq^Cs q-zS^C^Hq\ eq@-<S^ s<C- qbs zb sS eY<PbqCs ^Lq<Cq%PbeeS^Li

### 1.1 Motivation

r~zb^ - ^@3-qz^ 9: @q ... - e-q YCY4Cz..CC^ qS^HqC C^z YG q^S^L - ^@zPC...%S^S^L 4CS^Ls YG q^  
zPq-LP eCqS^z b4sCqf zb^ - ^@S^zCq <zb^ ..SZP zPCq s~q^@S^L C^fSj^ \ C^zi yPS b4sCqf zb^Q  
S^zCq <zb^ e-Sj beC^s zPC @bbq zb ^C...: %bHeb^@Cq zb^ ..PC^CfCq - ^C..<P- YC^LC- qS^S S^ ..PS^P  
e-sz C^eCq^<G < ^ 4C-eeY@ pC^HqC C^z YG q^S^L < ^ 4Cs-\ \ C@~e - s zPC eq<Css bHYG q^S^L  
<~s- Y% ^@zPC\ - eeS^L 4Cz..CC^ - <zb^s - ^@<b^sD ~C^<Si

, YPb-LP zPS <YS S <bqC>zPCq- qCsL^S < ^z @S Cq^<G 4Cz..CC^ Pb...- YfS^L 4CS^L..b-Y@YG q^  
- ^@Pb... \ -<PS^C..b-Y@ @ zPCs-\ G GSpz bH- YzPC-\ b-^z bHS^Hq\ - zb^ qd ~Sj@ Hq- \ -<PS^C  
zb \ -zP zPC...% P-\ - ^ eCqHq\ s b^ - seG-S <z-sWCj<CC@s zPC ^-\ 4Cq bHeCqCeZb^s - P-\ - ^ ~sG

HhQ Cqzq ebYzB^ - ^@C<SS^Q - W^L>\ - W^L zPCzq S^S^L bH ^%qS^HhQ^C C^z YC q^S^L \ b@CY- @z Q  
S^zC^sSfC eqp<Css fbH^C qI ~Sf^L \ S^B^s bHC^fS^p^\\ C^z szCesg..PS^P - Yb \ - W^S Sz - zS CQb^s-\ S^L  
eqp<Css zP-z qI ~Sf^S - YqC-\ b~^z bH^b\ e~z zB^~Yeb..Cq..PCq^s - P~\ - ^ < ^ YC q^ Hh\ - s\ - Y  
^~\ 4Cq bH^S^zCq <zB^s - ^@sPCq S^z~S^B^i

rCb^@bH Y^b s%szC <b\ Cs <YsCbz zPC- 4SS^%bHLC^Cq Y^zB^ zP-z- P~\ - ^ P-s>..SZPb-z P- f^S^L  
<P- ^LGS S eY^C C^zC@S^zB Ss szq<z-q>4C<- ~sC- Ys%szC s Y<WzPC<- e- 4SS^%zb LC^Cq Y^C <b^<Cezs - ^@  
zPC- 4SS^%zb YC q^ ^C..z-sW^S^ - s- \ eYCC <S^z\ - ^Cq 9:i

yPC Y<WbHs-\ eY^C <S^<%σ Yb P-\ eCq^ zPC..bqWbHqCsG qPCq^ - ^@\ - W^S Sz q zPCq @S <~Y zB  
<b^@<z C^eCq^ C^zs - ^@<bqCz \ Sz W^S>GseC-S W^Lsf^C Pb...@S <~Y Sz S zB^" ^@4-LS S^ qS^HhQ^C C^z  
YC q^S^L q^S^G qP - ^@<C zB \ - ^%eq4Y^ s S^PCq^z zB \ - P^S^C YC q^S^L 9:>s~<P - s q^@b\ sCC@S^L>  
szb~P- szS^%bq^s^ eY^δ C^z..bqW^PbS^G

X-sz%b...s \ eY^C <S^<%δ - %eqfS^C sbYzB^s zB zPC \ - W^S^%bH^z%eq4Y^ s zP-z- qC~s~ W^%  
zq zC@S^ zPS^ - qC >s~<P - s s^L ~YzB^s - ^@L- \ Csi R^ seS^CbHzP- z>zPCq^S^ szSY- sL^S^ < ^z L- e S^ Pb..  
zPS^ ..b~Y@ zq ^sYzC zB qC Yq.bq^@ eq4Y^ s s~<P - s S^@-szq^σ ~zB\ - zB^>PG YP< qC> q4bzS^s - ^@zPC  
q eS^%Lq^..S^L^ C@ bH ~zB^b\ b-s @f^S^L>..PS^P P-s 4C^ CszS^ - zC@zB P- fC- \ - qWz f Y-C bHbfCq |  
zqfS^B^ @bY^q^s 4%q C^E

R^ sPbq>zPC- 4SS^%zb sbYfCqS^HhQ^C C^z YC q^S^L eq4Y^ s S^ - sPbqz- \ b~^z bH^S G<b\ 4S^C@..SP  
Ss CfCq^S^ <q^S^L ^C@S^ qC Yq.bq^@ eq4Y^ s>s-LLGsz zP-z s- \ eY^C <S^<%εb~Y@- ssSz qS^HhQ^C C^z  
YC q^S^L S^ \ - W^L zPC ^C^z 4SL YC e Sz qI ~Sf^si

## 1.2 Objectives

yPCeq\ - q%Lb- YbHzPS^ ..bqW^S zB S^fCszL- zC zPC <- e- 4SS^S^s - ^@ebzC^zS YbH f- qCz%bHzq ^sHhQ^C Cq  
- qPS^Cz-zq^s HhqsbyS^L qS^HhQ^C C^z YC q^S^L eq4Y^ s ..SZP - Y^ S^C@^~\ 4Cq bHs-\ eY^S>- s..CY- s zB  
- ssC^s Pb..G <P - qPS^Cz-zq^ - ^@S eY^C C^z zB^ <- ^ S e- z eCqHhQ^C - ^<C S^ Y^ S^C@< z- C^fS^p^\\ C^zs  
s~<P - s zPC? C^e| S^@; b^zqYr~S^C 9: bqzPC, z- q^cC^W^9:i

[ bqC seC-S^ < W^%zPS^ ..bqW..SY- Yb C^eYbC zPC <- e- 4SS^S^s bH\ b@C^H^C qS^HhQ^C C^z YC q^S^L  
- qPS^Cz-zq^s ..SZPb-z eY^S^L - YbqzP^ s> ~zS^S^L - qPS^Cz-zq^s s^L S^f q zB zPbsC ~sC@ S^ YqC eqCQ  
zq S^C@ \ b@C^s zP-z- qC^" ^Cq^<C@ zB @b..^szq^ \ z-sW^ 9>: - ^@ Pb..@S Cq^z- zC^zB^ \ C^P- ^S^ \ s  
s~<P - s \ ~YSPC @- zC^zB^>~sC@ S^ zPCzq ^sHhQ^C Cq- qPS^Cz-zq^ 9:CE..SYeCqHhQ^C ..PC- eeY^C@ zB zPS^  
z^@C bHeq4Y^ i

R^ - @S^B^>q^S^C qP ..SY^4C <b^@<zC@ zB S^fCszL- zC zPC ebzC^zS YHhQ^C - @f- ^<C C^zs S^ zPC <- qf^z  
sz- zC bHzPC- q^i

yPC zq ^sHhQ^C Cq- qPS^Cz-zq^s zB 4C S^fCszL- zC@- Yb P- fC- ^ - ~zBqLq^ssSfC ^- z-q^ - s b^C bHzPCq  
sz- ^@< z H^ z-q^s> ..PCq^ <- ~s- Y^%S^s q^eqp@<C@> \ S^ S^W^L zPC <b^@S^B^s bHqS^HhQ^C C^z YC q^S^L  
eq4Y^ s - ^@eqfS^S^L - ...%zB S^ \ C@S^zC%C^e- ^@b^ <- qf^z ..bqW@b^C S^ zPC b^ S^Cq^ X s-4sCz bH

ep4Y\ s 9c:i y PCqC ..SY- Ysb 4C- zZ\ ezs zb qGc qP - ^@ @CfCbe b^S^C pX sbY-zb^s ..SP - zC^zb^  
\ b@Cxi y PS qGc qP P- s zPC- @@C@ @S <-Y%bHszSY4CS^L S^ - q zPCq~^Cf^YqC@- qG i

### 1.3 Contributions

y PS zPGsS HbYb..s qC^z <b^zqf4-zb^s 4%b PC^ Ci - Y9: zP- zS^zpb@- <Czq ^shbq\ Cq zb zPCqS^Hbq-C^ z  
Yc q^S^L s~4sCz bH@Cge Yc q^S^L ep4Y\ si y PC HbYb..S^L <b^zqf4-zb^s - qCeqcS^zC@=

- á R@C^zS <- zb^ bHsd ~C^<C\ b@CS^L 4bzzY^C^W^S^ zPC? C~Sb^ yq ^shbq\ Cq- qPSzCz-qC
- á Bf- Y- zb^ bH@S Gq^z zq ^shbq\ CqS eY\ C^z zb^s S^ b^ S^C qS^Hbq-C^ z Yc q^S^L
- á , eeY- zb^ bHzq ^shbq\ Cq b^ \ b@CqC b^S^C qS^Hbq-C^ z Yc q^S^L - YbqSP\ s
- á Bf- Y- zb^ bHsd ~C^<C\ b@CS^L S^ zq @S^L C^fSp^ \ C^zs

### 1.4 Structure of the Document

y PC" qz <P- ezCq bHzPS zPGsS S zPC S^zpb@- <zb^>..PCq zPC\ bzsf- zb^> b4UCzSfCs- ^@ zPCszq- <z~qC bH  
zPC @b<- \ C^z - qCeqcS^zC@

- y PCq\ - S^@Cq bHzPS @b<- \ C^z S szq- <z~qC@- s HbYb..s=
- ; P- ezCq | eqcS^zs - Yc qz z-qC qfSC.. b^ fcg zq ^shbq\ Cq- qPSzCz-qC> f| g qS^Hbq-C^ z Yc q^S^L  
\ CzPb@S - ^@f{g zPC sz zCqHqPCq q <b\ 4S^ zb^ bH4bzP - qG s - ^@- @S<-ssb^ b^ Ss ezbC^zS YS^ zPS  
4b@%bH..bqW
- ; P- ezCq{ CfeYS^s zPCep4Y\ zb 4Cz- <W@- ^@ zPCzPCbqCzS- YsbY-zb^ b^ fcg? C~Sb^ yq ^shbq\ Cq  
f- Y@ zb^ - ^@f| g b^S^C qS^Hbq-C^ z Yc q^S^L zq ^shbq\ CqS eY\ C^z zb^si
- ; P- ezCqJ eqcS^zs qS~ Ys b^ fcg 4 sC- qPSzCz-qC f- Y@ zb^> f| gS eY\ C^zC@- qPSzCz-q Y<P- ^LCS>  
f{g zq ^shbq\ Cq k QC q^S^L - ^@fJg zq ^shbq\ Cq ebY%bezS S- zb^i
- ; P- ezCqI @S<-ssCs zPC qS~ Ys zP- z..CqC b4z- S^C@ S^ zPS ..bqW- ^@ HqzPCq..bqWzb 4C @CfCbeC@
- ; P- ezCqv <b^<Y@Cs zPS ..bqW4%@Cs<q4S^L Ss \ - S^ - <PSCfC^ C^zs - ^@4%eqpebs^L Hz~qC S epfCQ  
\ C^zs zb <-qqC^z ..bqW



# 2

## Literature Review

yPS <P- ezCq ..SYeqsC^z - YCq z-q cCfSC.. b^ fcg zq ^shbq\ Cq - qPSzC-z-qS - ^@ f|g qS^Hbq\ C^z YC q^S'L \ CzPb@s - s ..CY- s - @CzCq\ S'- zS^ bHzPCsq S e- z- ^@ qCf ^<Cb^ Hz- qCpX qS C qPi yPCqC ..SY- Yb 4C- f{g bfCqfSC.. bHzPCsz- zCbHzPC- q- ^@ Pb...S ..SYs e- z zPCeqpebsC@ ..bqW

yq ^shbq\ Cq - qPSzC-z-qS ..SY4CzPC \ -S^ Hb<-s bHzPS 4b@%bH..bqW ^@ HbqzP- z qG sb^ - ^ CteYQ ^- zS^ bHPb...zPC%H^<zS^ ..SY4CeqsC^zC@ , Hbq~^@Cqsz ^@S'L zPCsq 4- sC- qPSzC-z-qC>- sCz bHz- zCbHzPC- qS eY\ C^z- zS^s ..SY4CeqsC^zC@ - <bq@S'L zb zPCsq @S Cq^z S eba- ^<CS^ zPC- qG bH@C@e YC q^S'L - ^@b^ zPS 4b@%bH..bqW

pCS^Hbq\ C^z YC q^S'L S zPC- qG ..PCqC zPC sz- @C@ - YbqfP\ s S^ zPS ..bqW- qC S^sCqC@ - ^@ HbqzP- z qG sb^ - 4qfHS^zq@<zS^ zb zPCs~4Cz S eqsC^zC@> 4Cbq\ bfS'L b^zb - ^ b' @bY-% YbqfP\ S kQC q^S'L - ^@- ^ b^@bY-% YbqfP\ S eqfS - YebY-%bezS S- zS^i

CS- Y%zPC- eeY- zS^ bHzq ^shbq\ Cq\ b@C S qS^Hbq\ C^z YC q^S'L ..SY4C- ^- Y%C@- s - ..PbY> S ..P- z ..SY4CzPCsz- zCbHzPC- q eqsC^z- zS^i

## 2.1 Transformer Architectures

$a^{\wedge}CbHPCSSs \sim Gs \dots SP; | | s \sim ^{\wedge} @bzPCqz\%Gs bH^{\wedge}C-q Y^{\wedge}Cz . bqW\%S zP-z zPC\%P- fC- " \ddagger C@^{\wedge} \sim \ 4CqbH^{\wedge}e \sim zs$   
 $\sim ^{\wedge} @b \sim ze \sim zs \sim s \dots CY-s - " \ddagger C@^{\wedge} \sim \ 4CqbH^{\wedge} - eeS^{\wedge}L szCesi yPS \setminus - V\%s \mathcal{S} @S \leftarrow \mathcal{Y} zb \text{ eq} \langle Css sD \sim C^{\wedge} \langle Gs \mathcal{Y}\mathcal{W}$   
 $sC^{\wedge} zC^{\wedge} \langle Gs \sim ^{\wedge} @zP-z \mathcal{S} \dots P\%o q \sim PSCz \sim q\%s \mathcal{Y}\mathcal{W} q \sim \sim \text{ff}^{\wedge} z^{\wedge} C-q Y^{\wedge}Cz . bqW\% 9 | : > Xry | s 9 \{ : - ^{\wedge} @Kp \} s 9 J :$   
 $P- fC 4CC^{\wedge} \sim sC@S^{\wedge} \text{ eq} 4Y \setminus s \mathcal{Y}\mathcal{W}^{\wedge} \sim z-q YY^{\wedge}L \sim LC \text{ eq} \langle Css S^{\wedge} Li$

$yq^{\wedge} sHq \setminus Cq \sim qC - \setminus b@CYzP-z \sim sGs sCYQzC^{\wedge} zB^{\wedge} > - \setminus C-P- ^{\wedge} S \setminus bH \langle bqCYzB^{\wedge} 4Cz . CC^{\wedge} \text{ ebs} \mathcal{S} B^{\wedge} s S^{\wedge}$   
 $\sim sD \sim C^{\wedge} \langle Gz \text{ } 4CzzCq \text{ } q\text{eq} \langle C^{\wedge} z \text{ } e \sim zCq^{\wedge} si \text{ } yPS \sim q \sim PSCz \sim qC \setminus - V\%s C \ddagger zC^{\wedge} sSfC \sim sC bHe-q YCY \text{ eq} \langle Css S^{\wedge} L >$   
 $qC@ \sim S^{\wedge} L zPC \langle b \rangle e Y \ddagger \mathcal{S} \% bH\mathcal{S} zq S^{\wedge} S^{\wedge} L zS C sL^{\wedge} S^{\wedge} \leftarrow ^{\wedge} z\%oy PC " qz zq^{\wedge} sHq \setminus Cq \setminus b@CY9 \langle \mathcal{E} \sim \langle PSC \rangle fC @sz zCQ$   
 $bHQPCQ q \text{ } q\%s \sim \mathcal{Y} S^{\wedge} zq^{\wedge} sYzB^{\wedge} z \sim sV \rangle \sim ^{\wedge} @sb \setminus CbH\mathcal{S} @Cqf \sim zB^{\wedge} s P- fCsPh . ^{\wedge} LqG z \text{ eq} \setminus \mathcal{S} CS^{\wedge} \sim qG s s \sim P$   
 $\sim s f\mathcal{S} B^{\wedge} 9 I : \sim ^{\wedge} @ f\mathcal{S} Cb \text{ eq} \langle Css S^{\wedge} L 9 v : i$

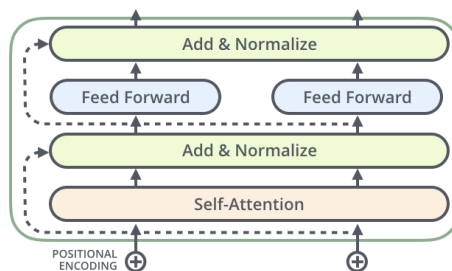
$yPCzq^{\wedge} sHq \setminus Cq \setminus - \%4Cse \mathcal{Y} S^{\wedge} zb z . b e \sim q\%s > zPC^{\wedge} \langle b@Cq \sim ^{\wedge} @ @C \sim b@Cq sz \langle W\%sPh . . ^{\wedge} S " L \sim q\%s | ic \sim ^{\wedge} @$   
 $| i | \text{ } q\%se C \sim SfC \% 3bzP \sim qC \setminus - @C \sim e bH\mathcal{S} CfCq YS^{\wedge} C^{\wedge} z\%s \sim YY \% Cpi yPC^{\wedge} \langle b@Cq sz \langle Wq \sim C\%SfC zPC S^{\wedge} e \sim z bHZPC$   
 $s\%s\mathcal{C} \setminus - Hq \mathcal{S} P \sim s 4CC^{\wedge} C \setminus 4C @ @ C @ S^{\wedge} - fC \sim zbq \dots SP zPC \text{ ebs} \mathcal{S} B^{\wedge} \sim YC^{\wedge} \langle b@S^{\wedge} L bHG \langle P e \sim q \text{ } bHZPC sD \sim C^{\wedge} \langle G$   
 $\dots PS \sim P \mathcal{S} \sim ^{\wedge} \sim q\%z P \sim z H\mathcal{Y} \mathcal{B} . s \sim e \sim zCq^{\wedge} zP \sim z \dots SY 4C YG q^{\wedge} C @ 4 \% zPC \setminus b@CY S^{\wedge} bq@Cq zb Wb \dots zPC \text{ ebs} \mathcal{S} B^{\wedge}$   
 $bHC \langle P bH\mathcal{S} s zbV\%s bq S^{\wedge} e \sim zsi$

$yPS \text{ } sC \sim zB^{\wedge} P \sim s zPC \text{ } e \sim q\%bsC bHC \ddagger eYS^{\wedge} S^{\wedge} L Ph \dots zPS \setminus b@CY 4C \sim \setminus C b^{\wedge} C bHZPC \setminus bsz \text{ } ebe \sim Yq @CCe$   
 $YG q^{\wedge} S^{\wedge} L \sim q \sim PSCz \sim q\%s > 4 \% @C \sim C \mathcal{S} S^{\wedge} L zPC \setminus bzsf \sim zB^{\wedge} 4CPS^{\wedge} @Ss @CfCbe \setminus C^{\wedge} z \sim ^{\wedge} @C^{\wedge} \sim \setminus Cq zS^{\wedge} L Ss 4C^{\wedge} C zsi$

### 2.1.1 Base Transformer

| icici,  $B^{\wedge} \langle b@Cq$

$yPC^{\wedge} \langle b@Cq Y \% Cq \text{ } eq \setminus - q\%Lb \sim Y\mathcal{S} zb C^{\wedge} \langle b@C \text{ } sb \sim qC sD \sim C^{\wedge} \langle Gs \sim ^{\wedge} @ \langle b^{\wedge} fCq zPC \setminus S^{\wedge} zb sz \sim zC fC \sim zbq \rangle$   
 $\sim \mathcal{Y} \mathcal{B} . S^{\wedge} L q\text{eq} \langle C^{\wedge} z \sim zB^{\wedge} s zb 4CYG q^{\wedge} C @ \sim ^{\wedge} @ @ z \mathcal{S} C^{\wedge} sB^{\wedge} s zb 4CqC @ \sim C @ Gs \sim qC | ic @ C e \mathcal{S} z s \sim ^{\wedge} C^{\wedge} \langle b@Cq$   
 $\sim s \sim ^{\wedge} C \ddagger \setminus eYI$



$GS \sim qC | ic = B^{\wedge} \langle b@Cq \sim q \sim PSCz \sim qC \text{ } f \text{ } p \text{ } C \text{ } e \text{ } q^{\wedge} zC @ Hb \setminus 9 u : g$

$B \sim \langle P C^{\wedge} \langle b@Cq S^{\wedge} zPC sz \langle W\mathcal{S} \setminus - @C \sim e bH \setminus . b s \sim 4Q \% Cpi R^{\wedge} zPC " qz b^{\wedge} C > zPC S^{\wedge} e \sim zs \dots SY e \sim ss zPq \langle LP \sim$   
 $sCYQzC^{\wedge} zB^{\wedge} Y \% Cq \dots PS \sim P \sim \mathcal{Y} \mathcal{B} . s zPC C^{\wedge} \langle b@Cq zb YbW \sim z bzPCq zbV\%s \dots PSC C^{\wedge} \langle b@S^{\wedge} L G \langle P se C \sim S \langle zbV\%s si$   
 $yPC b \sim ze \sim z bHZPS " qz s \sim 4Q \% Cq \dots SY 4C H @ \setminus zb \sim H @ @ Hq . : q @ ^{\wedge} C-q Y^{\wedge}Cz . bqWzP \sim \mathcal{S} zPC s \setminus C S^{\wedge} \sim Y$

ebsSb^si GbYb..S'L zP- z>S' ..P- z S Wb..^ -s- qSs@- Y<b^^C<zB^>zPCb-ze~z bHC <P s~4Q' %Gj ..SY4C  
 s~\ \ C@ ..SP Szs S'e~z zb - Yb...Hbq @CeCq^Cz .bqW ..SPb-z f- ^S'PS' L Lq @C^zs 9Di  
 rCYQzzC^zB^ P-s 4CC^ qC^ ^C@ fS - zC<P^S ~CWb..^ -s \ ~YSPC @-zzC^zB^>..PSP ..SYS\ eqfCzPC  
 \ b@CYS^ z.b ...: %s=

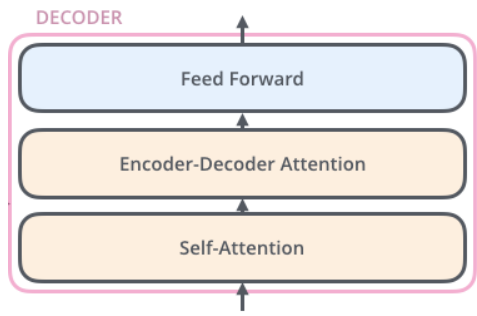
ci 3% Cte- ^@S^L zPC- 4SS% bHzPC \ b@Cyzb Hb<~s b^ @S Cq^z ebsSb^s

|i 3% eqfS^S'L ^~\ Cq~s qCeqS^z zB^ se- <Gs zb zPC- zC^zB^ Y%Gj G <P zbW^ ..SY4C qCeqS^zC@  
 S' e- q YCY\ - ^%zS C>- Yb..S'L \ ~YSeYC e- zCq^s zb 4C @S<bfCqC@

BssC^zS Y%a..P- zP- eeC^s S' "L~qC |ic ..SYP- eeC^ \ ~YSeYCzS Gs S' e- q YCY- ^@zPCqC..SY4C \ ~YSeYC  
 b-ze~zs <b\ S'L Hb\ zPC- zC^zB^ Y%Gj , YzPCb-ze~zs ..SY4C <b^<- zC^ zC@- ^@\ ~YSeYC@4%o ..CSPz  
 \ -zqf: zP- z ..SYsCqfC- s b-ze~zi

|icici3 ? C<b@Cq

yPC@C<b@Cq- qPSzCz-qC@GSL^ S'l ~ScsS S' qHq z-q'L- zPSqY%GjzP- zeCqHbq s \ ~YSPC @- zC^zB^  
 bfCqzPCb-ze~z bHzPC^ <b@Cqsz <W y PC- zC^zB^ S' zPC" qsz @C<b@Cqs~4Q' %Gj \ -%4C \ -sW@>HbqS'L  
 zPC- qf %zb @S\ Ss HqzPCq S'e~zs zP- ^ zPC <~qC^z b^G y PS\ \ -Wb sC^sC 4C- ~sC <bqCYzB^ sPb-Y@  
 b^%4C@b^C..SP eqfjqWb..Y@LC>Gsz 4YSPS'L - ^ S ebqz ^z qCYzB^sPSe bHk- ~s- Y%o  
 GSL~qC |i| @SeY%zPC @S Cq^<Gs Hb\ - ^ C^<b@Cq Y%Gj zb zPC @C<b@Cq Y%Gj



GSL~qC |i| = ? C<b@Cq- qPSzCz-qCfpCeqS^zC@Hb\ 9urG

yPS sz: <W<b^z S's - YbHzPC- zqf4-zGs bHzPC^ <b@Cqsz <Wb~P - s ebsSb^ - YC^<b@S'L - ^@ qSs@- Y  
 <b^^C<zB^si

|icici; , zC^zB^ [ C<P- ^S^s

yPC- zC^zB^ H^<zB^ S <b\ e~zC@4%ob^szq<S'L zPqC qCeqS^z zB^s bHzPC S'e~z zP- z- qC<- YC@zPC  
 k ~Cq%VC%o ^@, -YC fC<zqf Q>K>- ^@ Vg yPGsC- qC<qC zC@4%a ~YSeY%L zPC\ 4C@C@ S'e~z 4%o  
 zPqC \ -zqfGs zP- z- qCzq S^C@>Hbq C^ \ eYC zPqC Y^C q Y%Gpi



y PC-zzC^zsb^ H^<sb^ ..SY4C<-Y-YzC@S^ zPC HbYb..S^L ...%o..PCq d\_k S zPC @S C^sb^ bHzPC W%o  
fCz bq

$$Attention(Q; K; V) = \text{Softmax} \left( \frac{QK^T}{d_k} \right) V$$

, s eqfSb-sY%S^@S- zC@> \ ~YSPG @- zzC^zsb^ - Yb..s Hbq \ ~YSeY cCeCsC^z zsb^ s~4Qe- <Csi yPS S  
- <PScfC@4%zPC \ b@CYP- fS^L \ - ^%o zzC^zsb^ PG @s S^ G <P Y%Gq zb ..SYqCeCsC^z G <P s~4Qe- <G [ ~YSQ  
PG @- zzC^zsb^ S - sCz bH zzC^zsb^ <b \ e-z- zsb^s <b^<- zC^ - zC@- ^@ \ ~YSeY C@4%o ^bzPCq YG q^C@ ..CLPz  
\ - zcfi

[ ~YSPG @- zzC^zsb^ sCqfCs sCfCq Ye-qebsCs zPq-LPb-z zPC- qPSzCz-qC @CeC^@S^L ..PCq zPC Q>K>  
- ^@V fCz bq bqfLS^ zC- ^@ ..PCq zPC%o qCeP-CssC@ yPC bqfLS^ Y- qS^Y CqfC^ sz- zCs zP- z- zzC^zsb^ <- ^  
qCeCsC^z 4P- fSbq YW^s%z <S^ - ^@sC \ - ^zS^ szq-z-qCs S^ sC^zC^<Csi

### |icici? y PC H Y- qPSzCz-qC

y PCH^<sb^S^L bHzPCzq ^shbq\ Cq- qPSzCz-qC>- s eqCsC^zC@S^ " L~qC |i{>4LS^s ..SP zPC eqp<CssS^L bH  
zPC S^e-zsC \ -C^<CS^ zPC C^<b@Cq- ^@ @C<b@Cqsz <Wsi B- <P C^<b@Cqs b-ze-z S s-eeY C@- s zPC- zzC^zsb^  
fCz bq K - ^@V zP- z ..SY4C~sC@S^ C^<b@Cq@C<b@Cq- zzC^zsb^ zP- z ..SYPCé zPC @C<b@Cq Hb<-s b^ zPC  
- eeqepf zC qLsb^si

y PS eqp<Css S qCeG zC@ ~^zY- seC^S < s%o 4bYsL^ - Y zP- z zPC C^<b@Cq P- s <b \ eYzC@ Ss b-ze-z>  
- z ..PS^P ebsz S S e-ssC@ zb zPC 4bzbz \ @C<b@Cq ..PS^P - LLqL- zCs zPC qS- Ys - ^@ C 4C@S Ss b..^  
ebsSb^ - YC^<b@S^Li

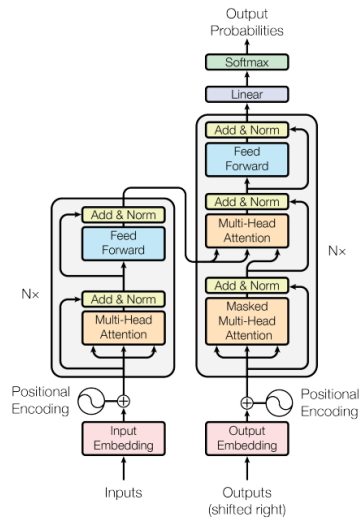
y PC s%zC \ <b^<Y@Cs 4%z- W^L S^ zPC- qf %b-ze-z bHzPC Ysz @C<b@Cq- ^@zq ^sYzS^L Sz b^C zbW^  
- z- zS CS^zb zPC @CsSfC@z%oC bH^Hbq \ - zsb^i y PS S - <b \ eYSPC@4%e- ssS^L zPC- qf %zPq-LP - H Y%o  
<b^<CzC@Y%Gq- ^@- r bH \ - † - <Sf- zsb^ H^<sb^ zP- z ..SY<b^ fCq C <P ebsS^Y zbW^ S^zb- eq4- 4SS%o  
- ^@sCC<S^L zPC b-ze-z 4- sC@ b^ zPC PSLPsz eq4- 4SS%o

### 2.1.2 Vision Transformer

Rz ...s ^bz Y^L 4C Hbq zPC ~^qG Y C@ ezbC^zS YbHzPCzq ^shbq\ Cq- qPSzCz-qC ...s szqCz-PC@S^zb bzPCq  
@Ce YG q^S^L " C@s s~<P- s <b \ e~zCq fS^S^i y PC, S^S^ yq ^shbq\ Cqf, S y g 9 I : eqfC@zb 4C- <b \ eCzsbq  
zb; ] ] - qPSzCz-qCs ..PSC~sS^L Yss <b \ e~z- zsb^ - YqSb-qCsi

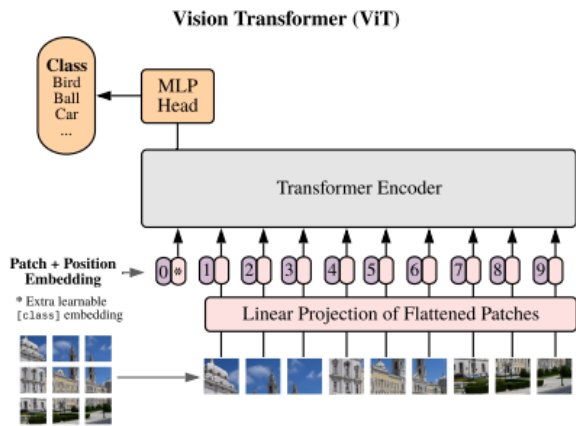
, \ - Ubqep4Y \ ..SP zq ^shbq\ Cq S^ fS^S^ S zPCY <WbH^@- <SfC4S sCs Hb-^@S^ bzPCq- qPSzCz-qCs>  
..PSP ..b- Y@H<S^z- zC Hbqzq S^S^L ..SP Yb... \ b-^zs bH@ z- i y PC~s- LCbH f- sz- \ b-^z bH@ z- zb zq S^  
zPC ^Cz .bqW^ ^@zPC^ " ^Cq- ^C Sz zb - s \ - YCqz sW S^L- zCs zPC- 4sC^<C bH b@CY- ss- \ ezb^si

, S y @S Cq Hb \ zPC bqfLS^ - Yzq ^shbq\ Cq @CSL^ S^ zP- z S^zC @bH .bq@s>zPC S^e-z S s - zzC^C@e- z-PCs  
bH ^ S \ - LCz bLCzPCq ..SP Ss ebsSb^ - YC 4C@S^Li ? CseS^C Ss Lbb@ qS- Ys ..SP YqL C- \ b-^zs bH@ z- >  
..bqW^ C@s zb 4C @b^C S^ bq@Cq zb qC@<C b fCq^ zS^L S^ s \ - YCq @ z- sCzi



$\mathcal{L} = \mathcal{L}_{enc} + \mathcal{L}_{dec} + \mathcal{L}_{attn}$

$\mathcal{L}_{enc} = \sum_{i=1}^N \mathcal{L}_{enc}^{(i)}$ ,  $\mathcal{L}_{dec} = \sum_{j=1}^N \mathcal{L}_{dec}^{(j)}$



$\mathcal{L} = \mathcal{L}_{cls} + \mathcal{L}_{reg}$

$\mathcal{L}_{cls} = \sum_{i=1}^N \mathcal{L}_{cls}^{(i)}$ ,  $\mathcal{L}_{reg} = \sum_{j=1}^N \mathcal{L}_{reg}^{(j)}$

$\mathcal{L}_{cls} = \sum_{i=1}^N \mathcal{L}_{cls}^{(i)}$ ,  $\mathcal{L}_{reg} = \sum_{j=1}^N \mathcal{L}_{reg}^{(j)}$

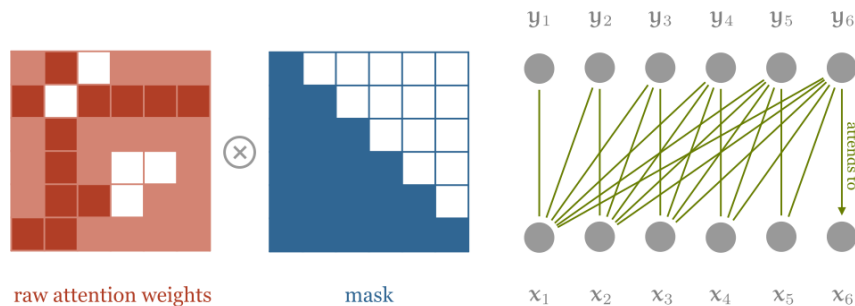
### 2.1.3 GPT-2

$y = \text{softmax}(W_{out} \cdot \text{encoder\_output} + b_{out})$   
 $z = \text{encoder\_output} \cdot W_{proj} + b_{proj}$   
 $z_{masked} = z \cdot \text{mask}$

$z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$   
 $z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$

$z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$   
 $z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$

$z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$   
 $z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$



$z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$

$z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$   
 $z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$

$z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$   
 $z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$

### 2.1.4 BERT

$y = \text{softmax}(W_{out} \cdot \text{encoder\_output} + b_{out})$   
 $z = \text{encoder\_output} \cdot W_{proj} + b_{proj}$   
 $z_{masked} = z \cdot \text{mask}$

$z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$   
 $z_{masked} = \text{encoder\_output} \cdot W_{proj} + b_{proj}$

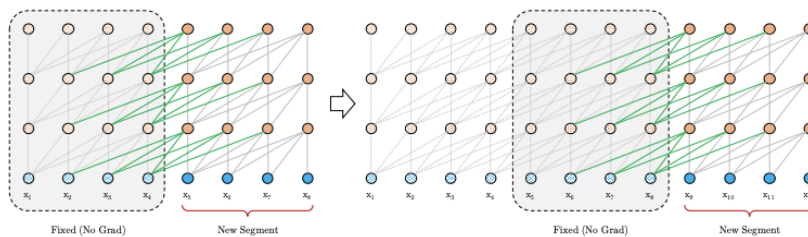
Pb... - <Gz- S' ..bq@<- ^ 4C- ^b-^ bq- fCq4 - <<bq@S'L zb Ss ebsSb^ b^ - sC^zC^<C bq - <<bq@S'L zb Ss s-qp-^@S'L <b^zCzi

### 2.1.5 Transformer XL

a Hc^>zPC Gsz- 4YsPC@ <b^zCiz Y^LzP bH e- qz~-Yq ecp4Y\ @bGs ^bz qz- S' zPC GssC^zS YWb..Y@LC zb sbYfCS>4G<- ~sC- e- qz~-Yq eSC-CbHf^Hbq\ - zb^ S s sz~-zC@zbb HqS' zPCe- sz zb 4Cq\ C\ 4Gq@>qS~Y S' ..P-z <- ^ 4C<- YC@- Y^L q ^LC@CeC^@C^<%ep4Y\ i

yb sbYfC zPC ecp4Y\ bH" iC@ <b^zCiz Y^LzP> zPC yq ^shbq\ CqQ X - qPSzCz-qC 9{: ~sGs qYzSfC ebsSb^ -YC^<b@S'ls S' bq@Cq zb - fbS@ ecp4Y\ s S' Pb... -zC^zb^ ..b~Y@ YbbW- z qCeG zS'L e- qz bH- sD ~C^<C- ^@szbqS eqfSb~s sD\ C^zs bH sD ~C^<C~sS'L qC~<qC^z ~^Ss zb 4C~sC@S' HqzPCqsD\ C^zs> \ -V^L zPC <b^zCiz Y^LzP ~<P Y^LCq

yPC- @f- ^z LCbHC eY%S'L \ C\ bq%S' zPCzq ^shbq\ Cq- qPSzCz-qC Ss qCeGsC^zC@S' "L-qC |ivi



$$GSL \sim qC | iv = fp Geq^zC@Hb \ 9{:g$$

## 2.2 Reinforcement Learning

yPC- LC^z - ^@zPC C^fSp^ C^z - qzPCz .b 4- sS< <b\ eb^C^zs bHp Xi yPC C^fSp^ C^z Ss zPCsz- zCq- <C S' ..PSP - ^ - LC^z YfCs - ^@- <zs b^ zPC @C~Sb^s S\ - Wsi , z G <P zS GszGe> zPC - LC^z b4sCqfCs zPC <-qC^z <b^@Sb^ bHzPC C^fSp^ C^z Ss S'P- 4Ss - ^@ zPC^ z Ws zPC - eepeq zC - <Sb^i ,, S'P CfCq% - <Sb^ zPC <-qC^z sz- zC bHzPC ..bq@ ..SY<P- ^LG R\ - %Pb..CfCq- Yb <P- ^LCb^ Ss b..^i

, HqC qC~SfS'L - <Sb^s> zPC C^fSp^ C^z sP- YqCz-q^ Ss <P- ^LGS S' zPC Hbq\ bH- ^C...sz- zC - ^@- qC..: q@sL^ - Y..PSP Ss - s<- YqzP- z Cf- Y- zCs l ~ ^zS- zSfC%Pb...Lbb@ 4CS'L S' - seG~S< sz- zC Ss Hbq- ^ - LC^zi yPSs Cf- Y- zb^ zCYs - ^ - LC^z ..PCzPCq Ss b^ zPC epeCq zq <Wzb - <PSCfC Ss Lb- Y ,, S'P zP- z s- S@- ^ - LC^zs e~qbsCSs zb \ - iS\ SCSs <- \ ~YzSfCqC..: q@bHC^ Wb..^ - s qz-q^i

3%LCzS'L - qC..: q@zPCqC Ss - ^ seG b^ ..PCzPCq bq^bz - ^ - <Sb^ Ss Lbb@ bq 4- @>4-z ^bz SHz Ss zPC 4Gsz - <Sb^i yPSs \ bzSf- zCs zPC ^CC@ Hbq CqeYq zb^ zb " ^@b-z ..P- z 4CP- fSbq sPb~Y@ 4C- @bezC@ yPSs <- ^ z W zPC sP- eC bHq ^@b\ fCqeYq zbqf\ b fCs bq \ b fCs zP- z - qC - YC @%qL- qC@ - s Lbb@ bezSb^s fLqCC@%b\ b fCs g

y PC b4CzsfC bHqS' HbqC\ C'z YC q'S'L Ss zb \ - VC zPC - LC'z YC q' Pb... zb - <PSCfC Ss Lb- Y R' zPS  
 sCzsb^>- ^ bfCqfSC... bHqS' HbqC\ C'z YC q'S'L - ^@sb\ C bHSs \ bsz qCCf ^z .. bqw' .. SY4C eqCsC'zC@

### 2.2.1 General Concepts

, sz-zC-z- e- qS~YqzS\ GszCe>st><b\ eYzC%@G<f4Cs Pb...zPC^fSp^ \ C'z Ss -z zPC<-qf^z zS\ G , Y  
 S' Hbq\ -zb^ S - Y.-%eqCsC'z S' zPCsz-zG Ob..CfCq zPC-LC'z \ -%b^Y%eCqCSfC e-qz bHSz \ -V'SL - ^  
 b4sCqf- zb^>ot>- e- qS Y@G<qezb^ bH sz-zC-zP-z \ -%bq \ -%b^z Y<V'S' Hbq\ -zb^i yPS \ G ^s zP-z  
 - ^ C'fSp^ \ C'z \ -%4CHY%bqe- qS Y%b4sCqfC@

B- <P C'fSp^ \ C'z P-s Ss b..^ sCz bHq-YCs> zP-z @CzCq\ S'C..PCzPCq- <zb^s - qC f- Y@ bq ^bzi ; PCss>  
 Hbq Cq- \ eY<-b^szq S's zPC \ bfC\ C'z bHG <P eSC-C - <bq@S'L zb sb\ C @Cqzsb^s - ^@ zb zPC ebsSsb^ bH  
 zPC qC\ - S'S'L eSC<Si yPC f- Y@ - <zb^s S' - <Cq- S' C'fSp^ \ C'z - qC <- YC@ zPC - <zb^Qe- <G

yq UCzbcfCs > f- Yb < YC@ qbYb-zsg - qC zPC sC \ C' <C bHSz-zCs - ^@ - <zb^s S' zPC .. bq@ zP-z Hbq  
 - ^ S'SS Ysz-zC> s0> \ -%4C qCeqCsC'zC@ 4% = (s0; a0; s1; a1; ...; sn)i yq UCzbcfCs \ -%4Cs- \ eY@ - ^@  
 sz-@C@ b^Y%os e- qS bH ^ .. PbYi

yPCq- qCz.b z%eCs bHpX, YlbqfP\ s>zPbsCzP-z - qC b^QbY%oa ^QbY%oa YlbqfP\ s  
 - qC zq S' C@ ~sS'L @ z- zP-z ...s <bYC-zC@ - <bq@S'L zb zPC \ bsz qC-C'z fCqSb^ bHzPC ebY%..PSC b' Q  
 ebY%oa YlbqfP\ s \ -%osC@ z- zP-z...s <bYC-zC@ 4%eqfSb-s S'z ^<Cs bHzPCebY%bq4%@ z- eqp@-C@  
 b-zs@C bHzPCebY%oa HHzPC @ z- S eqp@-C@ zbz- Y%b-zs@C bHzPCebY%zPC YC q'S'L eqp<Css Ss < YC@  
 a " S' CpX ..PS-P Ss zPC S' fCqC bHa ^S' CpXi

#### |i|ici, dbY%Cs

dbY%Cs - qC- \ - eeS'L Hb\ eCqCeZsb^s zb - <zb^s> zP-z sP- Y4Cz- V^ S' zPbsCsz-zCsi , ebY%e- ^ 4C  
 @CzCq\ S'SzS>..PC^ Sz b-ze-zs b^Y%oa ^ - <zb^ zb 4Cz- V^>S' ..PSP <-sCz < ^ 4C qCeqCsC'zC@ 4% (st)>  
 bq Sz < ^ 4C szb<P- szS' ..PC^ qCeqCsC'zS'L - eqp4- 4S%bHq G <P - <zb^i yPC ^bz-zb^>S' zPC Ysz <-sC  
 ..SY4C (jst)>..PCq (ajs) Ss zPCeqp4- 4S%bHzPC- <zb^ S' zS\ CqzCe t 4CS'L a HHzPCsz-zCs s>..PSP  
 < ^ 4C qCeqCsC'zC@ - s At = a>SHst = si

yPC qCeqCsC'z-zb^ bHbY%Cs ~sCs s-4s<qez >s^<Cz qCeqCsC'z zPCebY%e- q \ CzCq> \ G ^S'L zP-z  
 zPC b-ze-zs <b\ CHb\ - e- q \ CzCq @sb-qC> YV- ^G-q Y^Cz.bqWbq- zq ^sHbq\ Cq Hbq Cq- \ eYi

#### |i|ici3 pC.. q@ rL^ - Y

yPC qC.. q@ H^<zb^ @CeC'@s b^ zPC <-qf^z sz-zC bHzPC ..bq@> zPC - <zb^ z- V^>- ^@ zPC ^Cqz sz-zC bH  
 zPC ..bq@ yPCLb- YbH ^ - LC'z Ss zb \ - †S S C zPC <- \ ~YzsfC qC.. q@ bfCq- <Cq- S' zq UCzbcfCs

yPC qCz-q' H^<zb^>R( )>< ^ z- V^ \ - ^%Hbq\ s>..SP zPC \ - S' @S CqC'zS zbq 4CS'L ..PCzPCq bq ^bz  
 zPCqCs - eqp@%b^ S \ C@S zCqC.. q@ bq ^bzi yPC...%zb \ - zPC - zS- Y%qCeqCsC'z zPS Ss 4%e- \ \ S'L

zPC q... q@ S G <P zS CqzCe \ ~YSeY@ 4%o @S<b~^z H<z bq > zP-z @C<q sGs ..S P G <P zS CqzCe  
 q@~<S'L zPCS e- <z bHHz-q zS GzCes b^ zPC q... q@

$$R( ) = \sum_{t=0}^X t r_t: \quad f|i|g$$

„ PC^ 4G<b\ Gs c zPC q... q@ S ~^Sbq\ Hbq- YzS\ GzCes> - ^@ ..PC^ S S - ^%<f- YC 4Cz .CC^ CE- ^@  
 c<s~4s d ~C^z q... q@s ..SY@S\ S S P bfCq zS G

y b \ - †S S C zPC C eC zC @ qz-q^ S ..SY 4C ^C Gss- q%zb ^ @ zPC bezS - YebY%o i

**|i|ici; , - Y C G ^ <Sb^s**

, f- YCH^ <Sb^ S @S- zS zPC zbz Y- \ b~^z bHq... q@ zP-z S C eC zC @ zb 4C - <<- \ ~Y zC @ sz qS'L Hb\  
 - <Cz- S sz zC b q - s @C<qfC @ eqfSb~S%zPC C eC zC @ qz-q^i y PC \ - S b4C<zSfC bHqS Hbq<C C^z  
 YG q^S'L S bHC^ zPC GzS\ - zS^ bHfS H^ <Sb^> ..PSP S - zb-LP z sW

y PCq - qz . b \ - S z%Gs bHf- YCH^ <Sb^s> zPC" qz 4C S'L zPC b^ Q bY%<f- YCH^ <Sb^> V (s) > zP-z  
 b-ze-zs zPC C eC zC @ qz-q^ LsfC^ - <Cz- S sz qS'L sz zC - ^@ 4% HbYb..S'L zPC ebY% @ C S S^ si

y PCa ^ Q bY%o <Sb^Q - YCG- ^ <Sb^> Q (s; a) > GzS\ - zS zPC C eC zC @ qz-q^ LsfC^ - <Cz- S sz qS'L  
 sz zC > z V S'L - ^ - q4sq q%o <Sb^ - ^ @ zPC^ HbYb..S'L zPC ebY% HbqfCq

y PCq bezS - Y<b~^z Cq- qz > Q - ^ @ V > se C S% zP-z zPC%o q C HbYb..S'L - ^ bezS - YebY%o i

**|i|ici? [ b@Ck**

, \ b@CY\ S Ss Pb... - ^ C^ fSp^ \ C^z 4CP- fGs - ^ @ - Yb..s zPC- LC^z zb eqC@sz Hz-q sz zCs - ^ @ q... q@s  
 ..PSP <q zS zPC- 4S S%zb eY^ - PC @ [ b@Ck bHPC C^ fSp^ \ C^z @b ^bz - Y.: % C fSz sb - \ b@CY\ - %  
 4C YG q^C @ > 4-z S \ ~sz 4Cz- VC^ S zb <b^ s C q zS^ zP-z - YG q^C @ \ b@CY..SYP- fC- 4S s zP-z \ - % G @  
 zb s~4Q bezS - YeCq Hbq - ^ <C S - q C YC^ fSp^ \ C^z

„ bq S'L S - \ b@Cq C s Cz S'L S sS SY q zb ..bq S'L 4% zC Y- ^ @ Cq p q

**|i|iciB [ - qWf ? C S S^ d q b < C s s G s**

[ - qWf ? C S S^ d q b < C s s G s f [ ? ds g - q C ~sC @ zb q C e q s C^z q S^ Hbq<C C^z YG q^S'L e p 4Y\ s 4% @ C^ ^ S'L  
 zP-z zPC Hz-q sz zCs <b^ @ S S^ - Y% S @ C e C^ @ C^z b^ e- sz sz zCs > Lsf S'L zPC e q s C^z sz zC bq P [ s\_{t+1} j s\_t ] =  
 P [ s\_{t+1} j s\_1 ; ; ; ; s\_t ] i

[ ? ds z VC - @f ^z LC b HPC [ - qWf e p e C q %o ..PSP sz zCs zP-z zPC- LC^z P- s Wb..Y @ LC b H Y e- sz  
 S z C q <Sb^ - se C z s ..S P zPC C^ fSp^ \ C^z zP-z P- @ S • ~C^ <C S^ zPC Hz-q

y P S Hbq \ - Y S - zS^ S LsfC^ 4%o z-e YG=

á S QzPC s Cz b H e b s S Y C sz zCs

á A QzPC s Cz b H f Y @ - <Sb^s

á R QzPCq...q@H^<zb^

á P QzPCeq4 4SS%@Szp4-zb^ bHC <P zq ^szb^>qeqsC^zC@ eqf3-s%o s P(jst; at)

### 2.2.2 Model-Free Algorithms

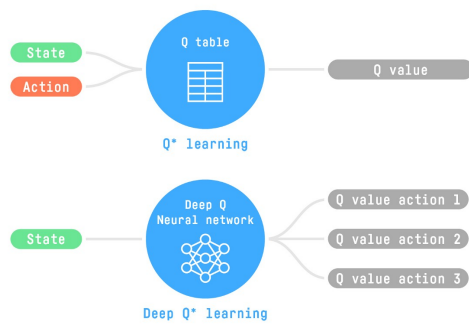
rb\ CbHzPC\ bsz <b\ \ b^ \ b@CqHC qS HqC\ C^z YC q^S'L - YbqzP\ s..SY4CeqCsC^zC@S^ zPs s-4sC<Q  
zb^>- Y^L ...SP zPCq qY zb^sPs zb qC~qq^z ^Cz.bqW - ^@ Pb...sC ~C^<C\ b@CS'L - ^@ zq ^shq\ Cq  
< ^ 4C~SS C@zb z <WCpX eq4Y\ si

|i|i,i, kQC q^S'L

? Ce k QC q^S'L 9J: bqlS^ zCs Hb\ sz ^@ q@k QC q^S'L - ^@ Ss b4CzsfCs zb \ - fS S C zPC bezS - Y  
b^QbY%o <zb^> qf YCH^<zb^> Q (s; a) > sz q^S'L Hb\ zPC <~qq^z sz zC>..PSP S - PSCfC@ 4%o ^@S'L - ^  
bezS - YebY%odfbq zb zPC ~sC bH^C-q Y^Cz.bqW - ^ P-sPQ 4C <b~Y@ 4C ~sC@S^ C^fSp^ C^zs ..PCqC  
zPC sz zCqe- <C- ^@- <zb^> qe- <C..CqCs\ - Y4%o eY\ C^z zb^ bHzPC k QH^<zb^>=

$$Q(s_t; a_t) = Q(s_t; a_t) + (r_{t+1} + \max_a Q(s_{t+1}; a) - Q(s_t; a_t)); \quad f|i}g$$

zP-z S <b^zS~ Y%oe@ zC@ ..b~Y@ 4C zPC YC q^S'L q zC- ^@ zPC @s~b^z H<zbj  
, <b\ e- qsb^ 4Cz..CC^ 4bzP \ CzPb@s S @SeY%o@S^ "L-q |iui



GS~qC |iu= ; b\ e- qsb^ 4Cz..CC^ k QC q^S'L - ^@? Ce k QC q^S'L fp CeqfzC@ Hb\ 9I:g

yPC GzS\ - zb^ bHsz zCQ <zb^ f-YCs - <bq@S'L zb - @sfCqC - \ b~^z bHsz zCQ <zb^ e- Sps S ~sC@S^  
@S<qCzC- <zb^> qe- <S- ^@zPC- <zb^ zP-z - PSCfCs zPCPSLPCsz k qf YC> argmax\_a Q (s\_t; a) >..SY4Cz- VC^i  
yPCzq S^S'L..SY4C@b^C 4%ob^zS~ Y%oe@ zS'L zPCk qf YC GzS\ - zb^s zP-z..SY4C<b\ Cb-q^C.Csz  
z qCz>yj - ^@- eeY%SL zPC \ G^ sl ~ q@ Cqpbq 4Cz..CC^ zPC <~qq^z k qf YC Hb\ zPC z qCz> LSfC^ 4%o  
zPCk H^<zb^ - 4bfC- ^@zPC eqf3-s k qf YC zP-z C^SzC@ 4C HqC zPC ^C...e@ zC>L = (y\_j - Q(a\_j; ))^2 >  
..PCqC - qC zPC ^Cz.bqW e- q \ CzCqi

? b~4Yk QC q^S'L 9v: - @qCsCs zPC eq4Y\ bHbfCqCsS\ - zb^ S^ sz ^@ q@k QC q^S'L..PSP qCs~ Ys  
S^ ~^sz 4C zq S^S'L - ^@- s-4bezS - YebY%oy Ps S - <b\ eY%PC@ 4%oS^S'L z.b ^Cz.bqW Hqk qf YC

Gsz\ - zsb^> zPC " qsz Hbq zq S'S'L zPC <-qq^z k Gf YGs> .. PS-P S zPC ebYS%o^Cz . bqW- ^@ zPC sCb^@ zb  
 Gsz\ - zC zPC k Gf YGs S' zPC ^Cfz sz- zC zb 4C ~sC@ S' zPC k QH^<zsb^> .. PS-P S ~e@ zC@ b^Y%o Hbq -  
 <Cqz S' ^-\ 4Cq bHszGes> .. SzP zPC bqlS'- Y^Cz . bqWb e- q \ CzCq> 4CS'L qCqCqC@- s zPCz- qCz ^Cz . bqW  
 ? Cc pC~qq^z k QCz . bqW 9u qCeY<C sb\ C bHzPCq ebYS%oY%Gp .. SzP qC~qq^z Y%Gp sS S' q zb  
 Kp} s 9J:~sS'L Ss PS@C^ sz- zCz szbqCS'Hbq\ - zsb^ 4Cz . CC^ zS GszGes- ^@zb YC q' e- zCq's zP- zC\ GqC  
 zPq-LPb-z- zq UC-zbq%o\ - eeS'L zPC k QH^<zsb^ .. SzP zPC- @sSb^ bH^ PS@C^ sz- zC b\ eb^C^zi y PSs  
 sPb-Y@^bz 4C b^HsC@ .. SzP sC ~C^<C\ b@CS'Li

**|i|i3 dqbS - YdbYS%a ezS S - zsb^**

dqbS - YdbYS%a ezS S - zsb^ fdda g 9D>s<<CC@s bzPCq, <zbq qfS\ CzPb@s YWC yq-sz pCLs^ dbYS%o  
 a ezS S - zsb^ fypdag 9\_: - ^@ , s%>Pqb^b-s - <zbq qfS f, |; g 9C> .. PS-P S' <bqbq zC 4bzP f- YCQ  
 bezS S - zsb^ YWC ? Cc k QCz . bqW> - ^@ ebYS%o bezS S - zsb^ 4%o-sS'L - ^ , <zbq zP- z b^zpbY Pb.. zPC  
 - LC^z 4CP- fCsi

y PSs z%oC bH- <zbq qfS - YbqP\ YC q's 4%o-sS'L Lq @C^z - sC^z zb \ - fS S C zPC - @f ^z LC> A\_t>  
 zP- z - <Cqz S' - zsb^> a\_t> P- s - L S' sz zPC q\ - S'S'L - zsb^s S' zP- z sz- zG y PSs S 4- Y^<C@ 4%zPC YLQ  
 eq4 4SS%bHz VSL zP- z - zsb^i - s LSfC^ 4%o

$$L( ) = E_t[\log (a_tj_s_t) \hat{A}_t]: \quad f|i|jg$$

y PSs YC @s zb ~^qCS 4CYC q'S'L> .. PCq zPC ~e@ zCS' zq S'S'L .. b-Y@CSzPCq 4CzbsYb .. bqbz ~^sz- 4Yi  
 yb Gsz- 4YSP - @LqC bH4 Y^<C> dda C\ eY%o- szq zL%Wb.. ^ - s zPC; YeeC@r~qpl- zC a 4UC-zSfC>  
 .. PS-P S Ss Pb.. YqC- ^ ~e@ zC zb zPC ^Cz . bqW sPb-Y@ 4G

CSzP> - q zsb S <- Y-YzC@ 4Cz . CC^ zPC ebYS%o zP- z P- s - YC @%ACC^ ~e@ zC@ - ^@ zPC ebYS%o zP- z  
 LC^q zC@ zPC @ z- zb l ~- ^z%zPC @LqC bH@SfCqL^<C 4Cz . CC^ zPCz . b ^Cz . bqW=

$$r_t( ) = \frac{(a_tj_s_t)}{old(a_tj_s_t)}: \quad f|i|jg$$

rS^<C zPS q zb b^sCq zPC eq4- 4SS% bHz VSL G <P - zsb^> Sz <- ^ qCeY<C zPC YLQ q4- 4SS%  
 bHC <P - zsb^> S' zPC Yss H^<zsb^i ; YeeS'L .. SYC^s- q zP- z zPC q zb .. SYsz- %ACz . CC^ - ^ S' zCqf- Y  
 [1 ; 1 + ] - ^@ zPC; YeeC@r~qpl- zC a 4UC-zSfC .. SYzPC^ 4CLSfC^ 4%o

$$L^{CLIP}( ) = E_t[\min(r_t( ) \hat{A}_t; clip(r_t( ); 1 ; 1 + ) A_t)]: \quad f|ivg$$

y PC\ S'S' ~\ 4Cz . CC^ zPC <YeeC@- ^@ ^b^q YeeC@ b4UC-zSfC Ss @b^CS' bq@Cq zb 4- Y^<C zPC ~e@ zC  
 - <bq@S'L zb Ss q zb - ^@ - @f ^z LC s~z- zsb^si y PC- @f ^z LC <- ^ 4C <- Y-YzC@ S' \ ~YSeY Hbq\ s YWC  
 KC^q Y, @f ^z LC BszS - zsb^ , ^@[ b^zCq - qb BszS - zsb^i



KL Cq Y-@f ^z LC SzS - zB^ fK, Bg <b\ e-zCs -@f ^z LCs -s

$$\hat{A}(s; a) = r + \gamma V(s_1) - V(s_0); \quad \text{f|iug}$$

..PSC [ b^zCQ - qb BszS - zB^ <b\ e-zCs -@f ^z LCs -s

$$\hat{A}(s; a) = rtg \quad V(s); \quad \text{f|iDg}$$

yPC " ^-YYBss H^<zB^ ..SY- Yb S<Y@C - , -YC Xbss H^<zB^ 4Gz..CC^ zPC f-YC eqpfS@C@ 4%zPC  
^Gz..bqW- ^@ zPC qG Yqz-q^ CfeCzC@ Hb\ zP-z sz-zC-s ..CY-s - ^ C^zpe%<bC <S^> r> zP-z PCes  
C^s-qC CfeYq zB^=

$$Loss = L^{CLIP}(\cdot) \quad k_1 L^{value} + k_2 \quad S(s_t); \quad \text{f|i_g}$$

pC~qq^<C<- ^ 4C - @@@ sS SY q%zb ..P-z P- eeC^C@ S ? Cce k Q Cz..bqW yPC f-YC - ^@ ebS%  
^Gz..bqW <- ^ sP- qC zPCsq 4b@%bq ^bz Qc:>..PSP \ - %4C ~sCHYzb ~^@Cpz- ^@ Hf z-qCs zP-z - qC qCf- ^z  
Hq 4bzP f-YC SzS - zB^ - ^@- <zB^ eqC@S-zB^i

## 2.3 State of the Art

[ bCq^@% q-PSz-z-qS YW zq ^shbq\ Cq P- fC ^bz P- @- 4SL S^zCpC-zB^ ..SP qS HbqC\ C^z YC q^S L  
eqp4Y\ s> \ - V^L S^ - ^ - qG beC^ Hbq CfeYq zB^i R^ zPs sC-zB^>sb\ C eS<Cs bH..bqWqCf- ^z zb zPC  
b4C<zSfC bHS fGszL-zS L Pb...^C...^C-q Y^Gz..bqW q-PSz-z-qS <- ^ S^~C^<C qS HbqC\ C^z YC q^S L ..SY  
4C eqCs^zC@

### 2.3.1 Stabilizing Transformers for Reinforcement Learning

yPCzq ^shbq\ Cqs <- e- 4SSCs zb eqp<Gss Y^L sC\ ~C^<Cs bHS Hbq\ - zB^ P- fC ^bz 4CC^ z W^ -@f ^z LC  
bHS^ qS HbqC\ C^z YC q^S L>@-C zb zPC @S <- YCs S^ bezS S- zB^ S^PCq^z zb zPs z^eC bH\ b@Cs> zP-z  
4G<b\ C C^fC^ \ bC^ ^bzSG 4YCS^ qS HbqC\ C^z YC q^S L @b\ - S^si

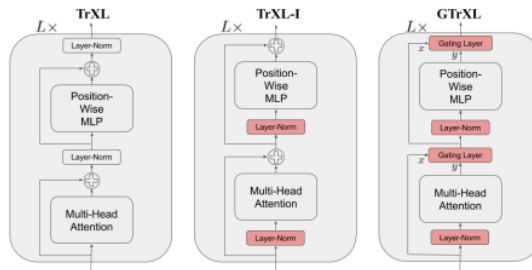
GL~qC |iD@CeSzs zPCK- zC@yq ^shbq\ Cq X- q-PSz-z-qCfKyqj XgQ |:>..PSP S^zC^@s zb sz- 4SS CzPC  
4- sCzq ^shbq\ Cq @GSL^ 4%eqpfS@S L L-zS L \ G-P- ^S\ s zP-z qeY<CzPC qS@- Y<b^ ^GzB^s - Y^Ls@C  
- sPSz S^ zPC ebszB^ bHzPC Y%q^ bq\ - YS- zB^ sC\ C^zsi yPs sCz bH- q-PSz-z-q YS\ eqpfC\ C^zs S  
zCzC@~sS L- , Q da - YbqzP\ S eY\ C^z-zB^ Q{i

Rz S 4CfC@zP-z zPs sbqz bHS^zS% - eeS L- Yb..s zPC\ b@CzB 4CS^S S @<YsCqzb - \ - qWfS ^  
ebS% ^@ zP-z S ..P%z zC^@s zb sz- 4SS C\ bC^G sS%oy PCL-zS L b-ze~z bH^ [ Xd bq sCfCzC^zB^  
Y%q S LsfC^ 4%

$$(W_g x) \quad x + (1 - (W_g)) \quad y; \quad \text{f|iCg}$$

..PCqy S zPC b~ze~z bHzPC Y%Gq - ^@ x S Ss S'e~z 4ChC Y%Gq bq) - YS - zSb^i

GS~qC |iD sPb..s Hb\ YH zb qLP> zPC ^C-Gss-q%eszGes zb zq ^sSb^ Hb\ zPC 4-sC zq ^shbq) Cq  
-qPSzCz-qS'zb zPC Kyqf X - qPSzCz-qG



GS~qC |iD= Kyqf X - qPSzCz-qCfp Geq'z@Hb\ 9 |:g

yPC Y%Gq ^bq) - YS - zSb^ ..SY4C ebsSb^C@ b^ zPC S'e~z szq \ bHzPC s~4\ b@ Ys f@C~b@Cq sz <Vg>  
-^@ s^<Cz..b YG q Y%Gq ...SY4C - eeY@ S' - sD ~^<C>- pCX) - <Sf- zSb^ ..SY4C - eeY@ 4ChC zPC  
L-zS'L \ G-P ^S\ si BfC^ ..SPb~z L-zS'L - eeY@> zPC sz 4SS%bHzPC \ b@CYS<qC sGsi

dY^z%bHL-zS'L S' eY C^z zSb^s ..GqCteGq C^zC@...SP>4-z- ^ - @ ez zSb^ bHzPCK-zC@pC~qq^z  
{ ^SfKp} g9J: - <PSCfC@zPC4Csz qS~Ys S' zPC? [ X-4Q Cb) - S' eqfS'L Pb.. \ ~<P 4CzCqzPS \ b@CYS  
<e-4SSCs bHq) C 4CqS'L Y^L q ^LC@GeC^@^<Ss - qG

yPS \ b@CYS - Yb qSsz ^z zb q ^b) C^fSp \ C^z sCC@S'L - ^@P%Gq- q \ CzCq bezS S - zSb^> \ - YS'L  
S zS G YHbq @G YS'L ..SP zPC ~^sz 4C b^@Sb^s bHqS'HqC C^z YG q^S'L C^fSp \ C^zsi

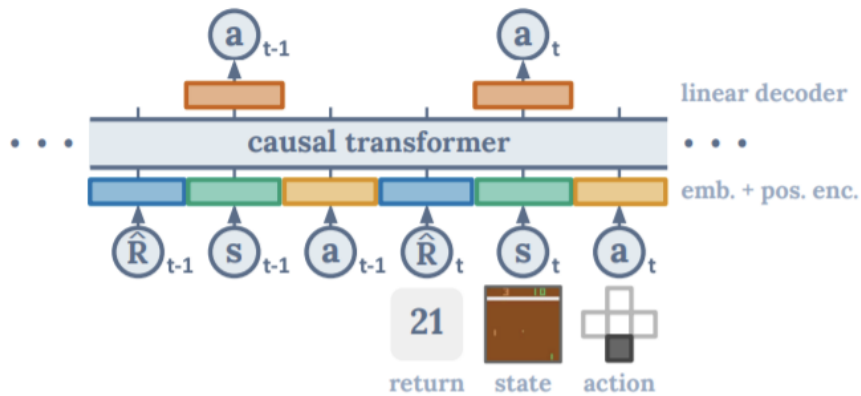
### 2.3.2 Reinforcement Learning via Sequence Modeling

rb) C YCq z-q z <VCS qS'HqC C^z YG q^S'L - s - s-eGqfS@ eq4Y\ 9J:> <b^fCqS'L qS'HqC C^z  
YG q^S'L eq4Y\ s S'zb s-eGqfS@ YG q^S'L eq4Y\ si yPC \ b@CY..SYz VC CteGq L \ Gs - ^@ q ^b)\  
L \ Gs sS Yq%zb ..P-z P-eeC^C@ S' zPC S' ebqz ^z , YP-Kb 9I: - ^@ ~sC zPC <- Y-YzC@ qz-q^s - s  
S'e~zs zb @S'L-Sp zPC | ~ Y%bHzPbsC sD ~^<G - ^@ zb 4C - 4C zb qCep@-C CteGzC@ qC.. q@si

[ b@Cq^@ %o qPSzCz-qS YVC zq ^shbq) Cq ..Gq - eeY@ zb zPS s-4QqG bHb" S'C qS'HqC C^z  
YG q^S'L 9>[v:>..PCq ^b sCY Y%P-eeC^s - ^@ zPC b^Y%YG q^S'L zb 4C @b^CS 4%o sS'L eqCsz 4YSPC@  
\ bfc) C^zs Hb\ - s-4Q bezS - Y@ z-4 sG

yPC \ - S' Hb~s <b^sSzs b^ C^ - 4S'L zq ^shbq) Cq zb YG q^ e-zzCq^s>-s zPC%eb sb ..CYS' | Xd>4%  
z YS'L - @f ^z LC bH~zbqLqssSfC \ b@CYS zP-z S' Sz-zC <~s Y%YVC zPC Kdy \ b@CY9c:> zP-z b^Y%  
YbW - z e-sz S'e~zs zb <b) e~zCzPC <-qq^z b^G> ~^YVC zPCL^Cq Y-sCbHzPC <b^zCz-z-s - ..PbY>-s ...s  
\ C^zS^C@S' sCzSb^ |ici

yPC - qPSzCz-q Y<PbsG - ^@ zPC S'e~z sD ~^<C - qC @GeSzc@S' "L-qC |i\_i



$$G_{\pi}^{\mathcal{E}}(s_t) = \mathbb{E}_{\pi} [ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t ]$$

$G_{\pi}^{\mathcal{E}}(s_t) = \mathbb{E}_{\pi} [ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t ]$

$G_{\pi}^{\mathcal{E}}(s_t) = \mathbb{E}_{\pi} [ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t ]$

$G_{\pi}^{\mathcal{E}}(s_t) = \mathbb{E}_{\pi} [ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t ]$

$G_{\pi}^{\mathcal{E}}(s_t) = \mathbb{E}_{\pi} [ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t ]$

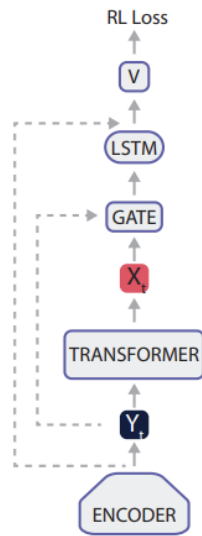
$G_{\pi}^{\mathcal{E}}(s_t) = \mathbb{E}_{\pi} [ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t ]$

### 2.3.3 Contrastive BERT for Reinforcement Learning

$G_{\pi}^{\mathcal{E}}(s_t) = \mathbb{E}_{\pi} [ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t ]$

$G_{\pi}^{\mathcal{E}}(s_t) = \mathbb{E}_{\pi} [ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t ]$

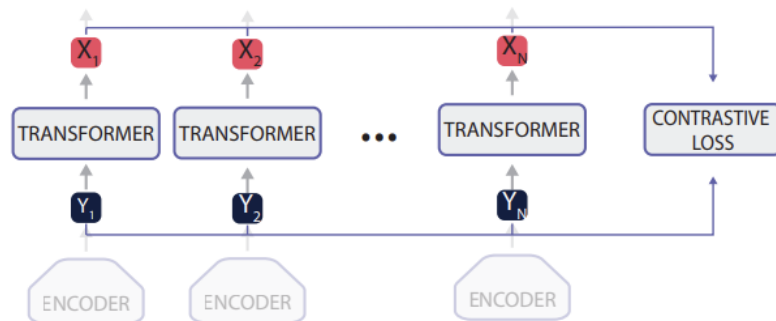
4S8q-zb^-YXry[ B^<b@SL \$ 4bqp..C@Hb\ 3Bpy 9: \ -VSL zPC\ b@CYH<-s b^ Ss s-qp-^@SL  
 zS, CszCes zb YC q^ -s sCC^ S' zPCHYY-qPSzCz-qCS' "L-qc|icE



GS~qC | icE ; b3BpX-qPSzCz-qCfpGeq^zC@Hb\ 9Dg

rCb^@> <b^zq szfC YC q^S'L \$ ~sC@Hbq zPC\ b@CYzb YC q^ LC^Cq YH z-qS - 4b-z Ss sz-zS ..SPb-z  
 @z- --L\ C^z-zb^>4C- ~sCzPC- ~zPbq qY%b^ zPCsD -C^zS Y^-z-q bHzPC@z- zb YC q^ zPCsCH z-qS  
 sS^<C@z- --L\ C^z-zb^ \$ ^bz zS CC <S^zi

yPC <b^zq szfC zq S'S'L \$ @b^C 4%λ -sVSL 15% bHzPC\ 4C@S'LS - ^@<b\ e-q^L zPC\ zb - sCz bH  
 ebsSfC - ^@ ^L- zfC C^-\ eYs Hb\ zPC @z- sCz zP-z-q e-ssC@ 4%zPC 4S8q-zb^-YC <b@Cq - ^@ zPC  
 zq ^shq\ Cq-s @SeY%C@S' "L-qc|icci



GS~qC | icc= ; b^zq szfC Yss S' ; b3BpX fpGeq^zC@Hb\ 9Dg

### 2.3.4 Online Decision Transformer

rS<C b" S' C qS' HhQ\ C' z YC q' S' L S Y Sx@ 4%zPC l -- Y%dbHzPC zq UC-zbqCs S' - @-z-sCz a ^S' C ? G<SS^ yq ^shhQ\ Cq 9\_>: Cte- ^@s b^ zPC qSc qP @b^ C 4%o PC^ Cz - Y 9: 4%o-^S%SL zPC eq<CssCs bH b" S' C qS' HhQ\ C' z YC q' S' L - s - eqCq S' S' L \ CzPb@ - ^@b^ S' C" ^Cq-^S' L S' - sS' LY<b^" L-q-zb^i , qeY%4-' Cq...SY4C - @@@ zb q<bq@ e-sz zq UC-zbqCs S' zPC C' fSp^ \ C' z ..PSP - qC - qC - sCz bH <b^sC~zSfC zq ^sSb^s> <b^szS~zC@ 4%zPC z-eYC (s\_t; a; RTG)> ..PSP ..b-Y@ 4C zPC eqfSb-s sz-zC> zPC ^Cz sz-zC> zPC - <Sb^ z-W^ - ^@zPC qz-q^ QbQbi a ^S' C" ^Cq-^S' L S - <<b\ eYSPC@ 4%ebe~YzS' L zPC qeY%4-' Cq..SP zPC zbe zq UC-zbqCs Hh\ zPC b" S' C @-z-sCz - ^@szbqS' L HhSP zq UC-zbqCs S' zPC qeY%4-' Cq~sS' L q ^@b\ %sCC-zC@ - <Sb^si yq S' S' L S - <<b\ eYSPC@ 4%e \ eY' L s~4Qq UC-zbqCs Hh\ zPC qeY%4-' Cq - ^@\ - †S S' S' L zPC YLQSWSPbb@bH zPC zq UC-zbqCs S' zPC zq S' S' L @-z-sCz ..PSC eq\ bzS' L CteYq zSb^ fS rP- ^^b^ C' zpe%o yPC b^S' C @G<SS^ zq ^shhQ\ Cq S \ bqC qS 4YS' \ bsz C' fSp^ \ C' zS S' ...s zCzC@b^> - ^@eqfS@Cs - sCz bH CL-zSfC qS~Ys zP-zPCe ~^@Cqz- ^@..P-z S - ^@S ^bz - ebssSfC - @f ^<CS' b^S' C CteYq zSb^ ~sS' L zq ^shhQ\ Cqi

# 3

## Methodology

Gbqsd ~C^<CLC^Gq zb^ S qS^Hbq\C C^z YC q^S^L>@C^b@Cq^Y%o q^PScz-z^Cs eqfCzPC\ sCfCs zb 4CzPC szp^LGSz ~sC <- sC bHzq ^shbq\ Gq> s^<C <- ~s- Y%@S-z zCs Pb...@C^SS^ \ - W^L S \ - @> 4- sC@ b^ zPC <- qf^z sz- zC bH \ b@CY- ^@ Ss sz- zC b^ - eqfS~s zS GszCei G-z- qC- <zb^s <b- Y@ b^%4C \ b@CC@ 4%o eqC@S^S^L G <P Hz- qCsz- zC- <<bq@S^L zb- eY^ bH <zb^si

? GseSczPCH<z zP- z k QC q^S^L - ^@ebY%bezS S- zb^ - eeq- <PGs - qC..S@C%o sC@> zPC%o qC~^@CqQ ~zSS C@..PC^ e- Sq@ ..SP sC ~C^<C \ b@CS^Li yPC \ - S^ <b^zq^1-zb^s bHzPS ..bqW- qC Y qL C%Hb<-sC@ b^ 4-S@S^L eqbH@H@b^<Cez - eeq- <PGs ..SP zPS Lb YS \ S^@..PSC C^s-q^S^L - @D -- zC f- Y@ zb^ - ^@ ebzC^zS YS eqfC\C C^zs zb ? C^SS^ yq ^shbq\ Gq

yPS sCzb^ S @fS@C@ 4Cz..CC^ zPCf^eCq\C C^zs <b^@-<zC@S^ zPCb^ S^CpXsCzS^L>zP- z C^e- ^@- ^@ f Y@ zC? C^SS^ yq ^shbq\ Cqs qS~Ys fcg- ^@zPbsCzP- z z- <WC C^eYq zb^ - ^@ebe~YqpX- YbqfP\ s f|g<b\ 4S^C@ ..SP sC ~C^<C \ b@CS^Li

### 3.1 Decision Transformer Validation

#### 3.1.1 Network

yPC4-sC^Cz.bqWzP-z...s~sC@zb qCep@-C? y s qS~Ys 4CLs's 4%qC-Gs's'L zPC^~\ CqS- Ybc@Cq bHZPC zS GzGes S' - sD ~C^<C-s ..CY-s - sCz bH <Sb^>sz-zCs>- ^@ qCz-q^sQbQb fpyK sg S' G <P zS GzGei yPC sz-zCs - ^@- <Sb^s - qC @q...^ Hb\ zPC @-z-sCz> ..PSC zPC pyK s - qC < Y~YzC@ 4%e~\ \ S'L zPC b4z-S'C@qC...q@ Hb\ zPC C^@ bHZPC zq UC-zbq%zb zPC <-qC^z ebS'zi , qCeqCsC^z-zSb^ bHZPS ^Cz.bqW <- ^ 4CsCC^ S' "L-qC |i\_i

yPCzS GzGes e-ss zPq-LP - ^ C 4C@S^L Y%Gq-zP-z sS eY%szbqS C 4C@S^Ls bH " †C@sS C^S' zPS <-sCzPC\ - †S ~\ ^~\ 4Cq bHS GzGes ebssSYC>eq@-S'L ^b Lq @C^zi yPCsz-zCs>- <Sb^s - ^@pyK s Lb zPq-LP Y^G q Y%Gq S' bq@Cq zb P-fCzPCs \ C@S C^sSb^- YZ%o

, HqLbS'L zPq-LP zPC Y^G q Y%Gq zPCsz-zC>- <Sb^ - ^@pyK sD ~C^<Cs - qC S^zCqG fC@S' zC\ ebq Y bq@Cq fGL st; rgt; at; st+ 1 iiii>- ^@zPC f-YCbHZPCzC\ ebq YC 4C@S^L-z zS GzGe t S - @@C@zb zPCe- qz bHZPCsD ~C^<CzP-z 4Cb^L zb zP-z zS GzGei X-%Gq bq\ - YS-zSb^ S - eeYC@4CbqCsC^@S'L zPCsD ~C^<C zb zPC @Cb@Cqsz-<W

B- <P @Cb@Cq S' zb\ ebsC@ 4%o \ ~YQP @-zC^Sb^ 4YbW ^@- \ ~YSQ^Gq eCq-Cezp^ ..SP KBX} -s Ss - <Sf-zSb^ H^<Sb^i yPC C^@ bHG <P bHZPbsCs-4QY%Gq S HbYb..C@ ~e ..SP qS@- YCb^<CzSb^s - ^@Y%Gq bq\ - YS-zSb^i GS'- Y%o Y^G q Y%Gq...SY<b^ fCq zPCb-ze~z @S C^sSb^- YZ%o S' zb zPC^~\ 4Cq bH ebssSYC <-Sb^s - ^@- ^ P%Gq 4bYQ- ^LC^z - <Sf-zSb^ H^<Sb^ ...s ~sC@zb YS S zPC - <Sb^s 4Cz..CC^ 1 - ^@1 S' - <bq@ ^<C..SP zPC C^fSp^ \ C^zs - <Sb^Qe- <G

, zC^zSb^ S \ -sV@ S' - ^ -~zbQLqssSfC \ - ^Cq zb C^s-qC zP-z G <P zS GzGe -zC^@s zb e-sz zS GzGes b^Y%o

yPC \ b@Cs ..Gq zq S'C@ Hb\ s-q z-P S' zPqC @Ss^<z sCC@s - ^@ ..SP zPC bqLS'- Ye-eCqs @CH-Y P%Gq-q \ CqCp> bHC^ %GYS'L - sYLPzY%PSLPq sz- ^@ q@ @CfS zSb^ 4Cz..CC^ q^s> 4-z - ^ D ~Sf- YC^z \ G^ eCqHq\ - ^<G

Bf-Y-zSb^ S @b^C...SP - " †C@pyK>zP-z S @Cq\ S'C@- <bq@S'L zb zPC\ - † qC...q@b4z-S'C@S' zPC C^fSp^ \ C^z \ ~YSeYC@ 4%o <b^sz ^zi yPCb^Y%oP- ^LCS' P%Gq-q \ CqCp ...s PCqC-sS^<CzPCeqCsC^zC@ qS~Ys S' sb\ C@z-sCzS ..Gq ^bz -PSCf- 4Y...SP zPC eqfS@C@z qCz pyK si

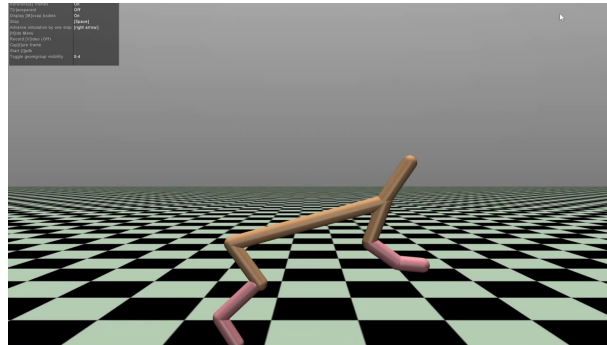
#### 3.1.2 Environments

ybeqpeGqY%f Y@ zCzPC-qPSzCz-q YS eY\ C^z-zSb^ - ^@zb Cq-e ^@~eb^ Sz>zPC~sC@C^fSp^ \ C^zs P-@ zb 4Csb\ CbHZPC~sC@S' zPCbqLS'- Ye-eCq GbqzP-z qG-sb^>[ ~Tb; bC^fSp^ \ C^zs JCE..Gq<PbsC^ @-C zb zPCq b4sCqf-zSb^ se-<C 4CS'L Y^G q ..PS&P \ - Ws zPC \ bqC C <S^z zb zq S^>-s sCC^ S' ; b4CqYQD - ^@Hq zPC @S <-Y%zP-z zPCq <b^zS' ~b-s - <Sb^Qe- <C- @s zb zPC D ~-zSb^i

yPC C^fSp^ \ C^zs ~sC@ ..Gq P- YQPCCz- P> PbeeCq - ^@ ... WCq zP-z Hq z-qC- sCz bHq 4bzS- sS Q

~Yzbc>..PbsC b4CzsfC \$ zb \ bfc Hsz - ^@S' - sz 4YC \ - ^Cq zb..: q@ - <Cq- S' Lb- Y yb @b zP- z zPC  
 \ b@CY^CCas zb b4sCqfC - sCz bH<bbq@S'- zCs>- ^LYCs> fCb<SSs - ^@ bzPCq s' bq@Cq zb @C<S@C Pb... \ ~<P  
 zbd ~Czb - ee%zb C <P bHzPC q4bz s PS^LGi , @CeSzSb^ bHb^C bHzPGsCs \ ~Yzbc \$ @CeSzC@S' " L~qC  
 {ic

y PC qC..: q@H^<Sb^s ecfS@C- ^ PSLPCq qC..: q@- <bq@S' L zb Pb... Hsz - ^@Pb... HqzPC q4bz - qC- 4YC  
 zb \ bfc zPC \ sCfCs S' zPC Y\ Sz z\ GszCes - Yb..C@



GSL~qC { ic= O- YQ PCz P B^fSp^ \ C'z

? - z sCz HbqzPGsCC^fSp^ \ C^zs ..CqCb4z S' C@Hb \ zPCe%zPb^ Y4q q%? JpX 9c:>- ^@ ..CqCbHzPqCC  
 z@Cs> \ C@S- \ > \ C@S- \ QqeY%o^@ \ C@S- \ QqeCqi

á yPC \ C@S- \ @ z sCz \$ LC^Cq zC@ 4%o ebY%zP- z - <PSCfCs b^C zPq@ bHzPC s<bqC bHzPC CqeCq  
 ebY%o

á yPC \ C@S- \ QqeCq @ z sCz \$ zPC<b^< zC^< zS^ bHzPC \ C@S- \ @ z sCz ...SZP - @ z sCz LC^Cq zC@  
 4%o ^ PSLPCq l ~ Y%ebY%o

á yPC \ C@S- \ QqeY% @ z sCz \$ LC^Cq zC@ Hb \ zPC qeY%Q- Cq bH- \ C@S- \ ebY%o \ -VSL Sz  
 \ bqC@SfCqCzP- z zPC q\ - S'S'L b^Gsi

y PC 4Gsz s<bqCs - qC CqeCzC@Hb \ zPC \ C@S- \ QqeCq @ z sCz >4~z sS^<C \ C@S- \ QqeY% @ z sCz - qC  
 zPC \ bsz @SfCqCzPGCs \$ sb \ CS^zCqsz S' <b \ e- q' L Pb... 4bzP bHzPGsC..b~Y@ \ b@CY- YCs zP- ^ bezS - Y  
 b4CzsfC> YW^bz LbS'L zb Hsz bqzbb sYb... Hbq Cq- \ eYi

### 3.1.3 Evaluation of Decision Transformer

{ici{i, [ G ^ pC..: q@, ^@ rz ^@ q@ ? CfS zSb^

yb ecfCq%cf- Y- zCb-qS eY\ C'z zSb^ bH C<SSb^ yq ^shh\ Cq- L- S'sz zPC 4- sCqCs- Ys> zPC \ b@CY... s  
 qC<b^szq<zC@Hb \ s<q zP - ^@ zq S' C@Hbq 100 zPb~s- ^@ szCes - <qss zPqC @S Cq^z sCC@si yC' q^s bH  
 zPC C'fSp^ \ C'z - qC~sC@ zb Cf- Y- zC zPC \ G ^ qC..: q@ bHzPC \ b@CY- ^@ Ss sz- ^@ q@ @CfS zSb^ - z CfCq%o  
 Cf- Y- zSb^ szCei



y PCsz ^@ q@ @CfS zsb^ \ Czqf - ^@ \ ~YseYC q^s - qC ~sC@sS^ <CzPC C^ fSip^ \ C^zs zb 4C sbYfC@P- fC  
 - @CLqCCbHszb<P- szSS%o ^@ ~^eqC@sSz 4SS%o ssb<S zC@ zb zPC \ >GseC<S Y%oS^ <CzPCqC\$ q ^@b \ S^SS YQ  
 S- zsb^ bHzPCsz qfS^L sz zGI

R^ zPSs sC<zsb^ qfS~ Ys ..GqC^ bq \ - YC@ zb 4Csl ~CC C@S^ 4Cz .CC^ ..P- z- q ^@b \ ebY%o.b-Y@ @b - ^@  
 ..P- z - ^ C^eCqz - LC^z ..b-Y@ @bi R^ zPC q \ - S^S^L sC<zsb^s - 4sbYzC qC..: q@s ..GqC eqfS^zC@ S^ bq@Cq zb  
 Csz 4YSP - \ bqC @Cqz <b \ e- qfb^ 4Cz .CC^ Pb...- qPSzCz-q Y<P- ^LGs - ' Cz zPC \ b@CY y PC^ bq \ - YC@  
 qfS~ Ys - qCb4z- S^C@ S^ zPC HbYb..S^L HsPS^=

$$] \text{ bq} \setminus \text{ YC@pC} \cdot \text{ q@} = \frac{\text{pC} \cdot \text{ q@} \quad [ \text{ S}^S \setminus \text{ pC} \cdot \text{ q@} ]}{\text{B} \ddagger \text{ eCqz pC} \cdot \text{ q@} \quad [ \text{ S}^S \setminus \text{ pC} \cdot \text{ q@} ]} 100: \quad f\{icg$$

**{ici{i3 ; b^zCiz XC^LzP**

? y - Yb eqfS^zS - ^ - ^- Y%ss bHPb...<b^zCiz YC^LzP S \ e- <zs zPC eqf4Y \ s S sbYfCs S^ zPC- z qfC^ fSip^ Q  
 \ C^zs>Pb..CfCqzPbsC C^ fSip^ \ C^zs - qCzq S^C@4%e \ eY^L1% bHzPC? k ] pCeY%? - z sCz 9 ] :>..PS^P  
 S^bz - ^ C^eCqssSfC - \ b^z \ G ^S^L zP- z zPCb4z- S^C@ qfS~ Ys \ - %bq \ - %dbz 4C C^eY^S^C@4%e \ eY^L  
 Y<W

GbqzP- z qC sb^ zPS ..bqW@sS^L-S^PGs b^ Pb...<b^zCiz YC^LzP ..b-Y@ - ' Cz @S GqC^z sbqzs bH@ z- sCzs  
 - <<bq@S^L zb zPCq ] -- Y%o ^@ @SfCqS%b^ <C- L- S^ eqfS^L zPC - sCHY^Css bHzPC <PbsC^ C^ fSip^ \ C^zsi

[ b@CfS ..GqC zq S^C@ C^ - <zY%Pb... zPC%a.b-Y@ S^ zPC bqfLS^ - Y <b^ @S^s^s> P- fS^L b^ Y%azPqCq <b^zCiz  
 q \ bfC@ ..PS^P \ - Ws zPC \ b@CYszbe \ b@CfS^L zq UC-zbqfS - ^@ sz qz YbW^L - z G <P zS CszGe b^ - ^  
 S^ @SfS@ - YHsPS^i

**{ici{i; y- qLCz [ b@CfS^L**

GbqzPCqG sb^s eqfS^zC@S^ zPCeqfS^s s-4sC<zsb^>- ^ Cf- Y- zsb^ bHPb...? y - <PSCfCs - @C^ @b4UC-zSfC  
 zP- z ...s <PbsC^ 4%Szs ~sCq ..SY4C <b^ @- zC@ S^ - \ bqC C^zC^sSfC HsPS^ zP- ^ ..P- z ...s @b^C S^ zPC  
 bqfLS^ - Y- qS^Yi yPS eqfS^C@ - sz~@%b^ Pb...@ z- S \ e- <zs b^ S^CpX - ^@ <b^ <Y@Gs b^ ..PCzPCq bq^bz  
 zPC S^zqb@<zsb^ bHC^eYq zsb^ S^ ^G<Css- qfS^ sD ~C^ <C \ b@CfS^Li

„ PC^ Cf- Y- zS^L zPC \ b@CfS - sCz bHD ~- Y%se- <C@ qC..: q@ b4UC-zSfCs ..GqC <PbsC^ - ^@ zPC - fCq LC  
 qC..: q@ b4z- S^C@ - <qss zC^ C^ fSip^ \ C^z Cf- Y- zsb^s ...s <b \ e- qC@ S^ bq@Cq zb ~^@Cqz ^@ SHzPC \ b@CY  
 <b-Y@ ~^@Cqz ^@ Pb...zb HbYb... seC<S^ <b4UC-zSfCbqS^z <b-Y@ CfC^ C^zq ebYzC - ^@s-qe- ss Szs b4UC-zSfCsi

y PC@S GqC^ <C4Cz .CC^ zPCb4z- S^C@ qC..: q@ - ^@ zPCz- qLCz qC..: q@ S^ ~sC@ zb b4z- S^ <b^ <Y^sb^s - 4b-z  
 Pb...G <P z@C bH@ z- sCz eCfHq ] s - <qss C^ fSip^ \ C^zsi

### 3.1.4 Proposed Architectural Changes

{iciJi, dbssb^ - YB \ 4C@S'Ls

, seqfS~s%CfeYS^C@>S\ GzGes..Gc^~\ 4Cq@- <bq@S'L zb zPCq ebsSb^ S' - sD ~C^<C- ^@CfeYS~S%  
 C^<b@C@- <bq@S'L zb zPC \ - †S \ ~ @S\ C^sS^ bHzPCzq UCzbcq%or S^<CG <P z-eYC bHzPqCC zbW^s fsz- zC>  
 - <Sb^>pyKg..b-Y@P- fCzPCs- \ CC\ 4C@S'L>zPS- eech- <P ..b-Y@H<Sfz- zCzPCYC q^S'L bHzPC qY zSfC  
 ebsSb^S'L 4Cz..CC^ z-eYsi , " †C@ \ - e ..b-Y@ ^bz b^Y%Y<W• C†S4Sfz%4Cz..CC^ zS\ GzGes> 4-z ..b-Y@  
 ^CC@- CfeYS~S\ S\ S\ HbqzPC C^fSp^ \ C^z zb 4C @CzG\ S^C@

R^ qG YqC ecp4Y\ s zPCqC ..b-Y@ 4C ^b ...%zb GzS\ - zC Pb... \ - ^%zS\ GzGes ..b-Y@ 4C ^C<Gss- q%zb  
 sbYfC- ecp4Y\ - ^@CfC^ SHzPCq...s>LC^Gq YS- zb^ 4Cz..CC^ ecp4Y\ s ..b-Y@ 4C sCfCqC%P~qz 4%zPSi

R^ bq@Cq zb sbYfCzPS ecp4Y\ >s^ - sbS@ YebSb^ - YC\ 4C@S'Ls ..GcS^zcp@- <C@S' - YC q^ - 4YCHsPb^>  
 P- fS'L zPC HbYb..S'L < Y~Y zb^ HbqzPCq qCeqsC^z zb^ se- <C 4CbqCzq S^S'L=

$$p_t = \sin \frac{t}{10000 \frac{2k}{d}} ; \quad f\{i|g$$

$$p_{t+1} = \cos \frac{t}{10000 \frac{2k}{d}} ; \quad f\{i|g$$

..PGCk S zPCqY zSfC ebsSb^ S' zPCPS@C^ @S\ C^sS^ bHzPC- qf %zb 4C ecp@- <C@- ^@d S zPC \ - †S \ ~\  
 @S\ C^sS^ bHs- S@- qf %

y PC~s-- Yq sb^ zb ^bz zq S^ ebsSb^ - YC\ 4C@S'Ls S zP- z YzCq ebsSb^s ..b-Y@ 4C<b\ C- ^@Cq q S^C@  
 @-C zb 4CS'L -sC@Yssi R^ zPS < sC-sS^<CzPC@ z- sCz <b^z- S's q-^s Hb\ sz- qz zb " ^SP>S 4C<b\ Gs G sCq  
 zb e-ss zPq-LP - YebSb^s bH zq UCzbcq%S^ - ^ ~^SHq\ HsPb^ \ - W^L YC q^ - 4YC\ 4C@S'Ls zPC 4Gz  
 - eech- <Pi

a ^CbHzPC \ - S' - @f- ^z- LGs bHzPS z@CbHC^ <b@S'L S zP- z zPC@Sz- ^<C 4Cz..CC^ zS\ GzGes S q@- <C@  
 ..SP zS\ G \ - qW^L - ^ S\ ebqz- ^z @SzS^<S^ 4Cz..CC^ zPC C^@- ^@zPC 4LS^S'L bH- ^ CeSb@C> ..PGC  
 -zz- S^S'L G q%sz- 4Sfz%S <q- <S Y..SP - " ^SPS'L sz- zC...PGC...cp^L - <Sb^s - qYss e~^SPS'L - ^@\ bqC  
 YS S^C@

{iciJi3 K- zS'L

y PC bS'L - Y? y - qPSzCz- q HbYb..s G <P bHs s- 4Q%Gp 4%o qS@- Y<b^GzSb^ zP- z s- \ s Ss S e-z  
 ..SP zPCs- 4QY<Wb- ze-zi yPS P- eec^s z..SCS^ G <P @C<b@Cq=

$$out_1 = ] bq\ Att(input) + input ; \quad f\{iJg$$

$$out_2 = ] bq\ MLP(out_1) + out_1 ; \quad f\{iJg$$

yPCz.bS eY\ C^zC@<P- ^LGs sz- q 4% eeYSL zPC ^bq\ - Y- zS^ Y%Gf 4CHqCC <P s-4Q%Gf

$$out_1 = Att ] bq\ (input) + input; \quad f\{ivg$$

$$out_1 = MLP ] bq\ (out_1) + out_1; \quad f\{iug$$

Rz S P%bzPGSS C@zP- z..PC^ Y q^S^L>zPCeqp@-<C@- <zS^ Hbq- sCz bHsz- zCs> (js<sub>t</sub>; :::; s<sub>1</sub>)>- eeqf†S  
 \ -zCs zPC- <zS^ zP- z..b~Y@4Ceqp@-<C@4%SS eY%oS^L zPCYsz sz- zC> (js<sub>t</sub>)>- eeqf†S -zS^L zPCebY%  
 zb b^C zP- z..b~Y@- <z- <bq@S^L zb - [ - qWf ? C<SS^ dqb<Gss>S\ eqpfS^L zPC<- e- 4SS%bH - \ b@CYzb  
 Y q^ ...SPb-z bJcGqCYS^L b^ - ^ C^zC CeSb@C bq- zq UC<zbq%z- VC^ Hb\ b^G yPS\ - Ws zPC\ b@CY  
 <- e- 4YbHqC <S^L 4CzCqzb zPC~^C^eCzC@ yPS P- eeC^s>4C<- ~sCzPCS^SS YC^<b@S^L S^ ^bzS e- <zC@  
 -s\ ~<P 4%δbq\ - Y- zS^ beCq zS^si

rS^<C zPC 4S - @f- ^z LC bHzq ^shbq\ Cq S ^bz zPCq - 4SS%zb qC < 4-z zPCq - 4SS%zb \ b@CYsCQ  
 l~C^<G>- ^ S eqpfc C^z zb Ss \ C bq%o. SY4C- @C@4%qCeY<S^L zPC qSS@- Y<b^<GzS^s ..SP L- zS^L  
 Y%Gf>zP- z- qLSfC^ 4%

$$(W_g x) \quad x + (1 - (W_g)) \quad y; \quad f\{iDg$$

..PGC %^ @ † @C^bzC zPC b-ze-z- ^@S^e-z bHzPC s-4Q%Gf..PGC L- zS^L S - eeY@zbi , G sCq fS-- Y  
 S- zS^ bHzPS <P- ^LGs S eqsC^zC@S " L~qC | iD

BfC^ zPb-LP zPS <P- ^LC- S s zbS eqpfc \ C bq%o ^@- Yb..Hbq- \ bqCsz- 4Y...%zb \ b@CYsD ~C^<G>  
 zPC zq @CQ S - ^ S^<qC sC bHc\_DcJJ e- q \ CzCq b^ - zPqCCY%q zq ^shbq\ Cq ..PS P ..SYsCfCqC%P~q  
 Ss zq S^S^L zS G

, ^bzPGq <P- ^LC zP- z- S s zb S^<qC sCsz- 4SS%S zPCS^SS Y- zS^ bHzPC 4S sCs zb - f- YCs-eCfqbz  
 <Cp>S^ zPC<- sC bHzPC <~qC^z C^eCf\ C^zs>Cfi

## 3.2 Online Transformer Evaluation

„ SP zPC b4CzsfC bHcf- Y- zS^L Pb..b^Y^CpX- YbqfP\ s eCqHq\ C@S^ sD ~C^<C \ b@CS^L>zPS sCzS^  
 eqsC^zs - \ CzPb@bYl%b^ Pb...zb - @ ez zPbCs \ C- YbqfP\ s S zb zPC sD ~C^<C \ b@CS^L sCzS^ Li

### 3.2.1 Environments

yPC <PbSC bHC^ fSb\ ^ C^zs S^ zPS sCzS^ @S Cq Hb\ zPC eqfS~s sCzS^> @-C zb zPC @S Cq^<C S^  
 b4CzsfG yPC- YbqfP\ s zb 4C z- <W@- qC ~sC@S^ C^ fSb\ ^ C^zs ..SP @S<qCzC - <zS^ se- <G>- ^@ Hbq  
 zP- zq sb^ zPC<PbS^CS^ C^ fSb\ ^ C^zs ^C@C@zb eqpfS^C- ^ b4sCqf- zS^ se- <CzP- z..b~Y^ ^bz qd ~SfCPSL P

<b\ e~z- zsb^- Yqsb-qcs zb 4C Sgq zC@ zPq-LPi ? S<qzC - <zsb^ se- <G ..SP @S Cq^z%P- eC@ qC.. q@ H^<zsb^s - qC- 4Y zb - ssCs Pb...zPC \ b@CY..b-Y@ 4CP- fCS @S Cq^z <b^@zsb^si

; - qzdbYCS - s C^fSip^ C^z..PbsCb4UC-zSfCS zb 4- Y^<C- ^ S^fCqC@ebYCS - \ bfS^L< qfi , fbS^S^L Ss HYHhq- ^bzPCq zS CszCe eqfS@G - s<- Yq qC..: q@ bHc>..PS-P S Y S@ 4%zPC \ - fS ~\ ^~\ 4Cq bH zS CszCes S^ zPC C^fSip^ \ C^z>ICEE

HH ^ - fCq LCqC..: q@ bHc\_I S b4z- S^C@S^ cEECeSb@Cs zPC C^fSip^ \ C^z S <b^sS@Cq@sbYfC@>Pb..CfCq ..C ~sC cEECeSb@Cs Hhq Cf Y- zsb^ sS^<C zPC C^fSip^ \ C^z S ^bz zbb <b\ eYq zb ...: qf ^z Cf Y- zS^L S^ cEE@S Cq^z CeSb@Cs eCq zq S^S^L ~e@ zG yPC \ S^S ~\ s<bqC S G sS%~q- ssC@ 4%zPC S eY C^zC@ \ b@Csi

[ S^S^Lq@ <b^sSzs bH-s^L - ^ N N qb\ - ^@ " ^@S^L zPC \ - qV@ C^Si R^ zPS <- sCzPC sz- q- ^@ C^S ebsSb^ - qC " fC@ - ^@ zPC b4sCqf- zsb^ se- <CS Y^G fC@ zb - Yb...Hhq z\ ebq YC <S^<%opC..: q@s - qC^b...se- qC- ^@ b^%LSfC^ HzPC- LC^z qC <PG Ss Lb- Y R^ bq@Cq zb s-<<C@>zPC- LC^z ^CC@s zb eY%o eCqf<z%bq Sz ..SYsS eY%L- S^ - qC..: q@ bHc Cq - z zPC C^@ bHzPC CeSb@C> YC q^S^L ^bzPS^Li , Y S@ ^~\ 4Cq bHzS CszCes S ~sC@S^ bq@Cq Hhq zPC C^fSip^ \ C^z zb YC q^ - ^ C <S^z \ CzPb@ zb Ss Lb- Y

, ^%aq @S^L S eY C^zs - szb^W - qVz zq @S^L C^fSip^ \ C^z zP- z ~sG - @- z- sCz bH <Cq- S^ szb^W - ^@ qCep@-<G Ss f- f zsb^ - <fss zS G>..PSC ~s^L - sCz bH^ ^- ^<S YS^@S- zbq zb qCeCsC^z zPC <-qf^z sz- zC bHzPC C^fSip^ \ C^z - Y^L ..SP zPC <-qf^z szb^WefC G yPC ~sC@ @- z- sCz ~sG - ^ Cf Y- zsb^ bHzPC szb^WefC bHK - \ Gzbe @ S%~ <fss z..b %G q@ -C zb zPC ePC^b\ C^b^ zP- z zPS szb^W4C<- \ G

, z CfCq%sz- zCzPC \ b@CYP- s zPC bezsb^ zb CszPCq

á rPbz - ebsSb^>SHz ...: ^zs zb sCYsP- qS ..SP PSLP f- YCzP- z S 4CXCfG <- ^ 4Ce~qP- sC@ YzCq- z Yb..Cq f- YG

á Xb^L - ebsSb^>SHz ...: ^zs zb 4-% ebsSb^ ..SP Yb...f- YCS^ zPC C^eCz- zsb^ zP- z Ss f- YC ..SY S^<qC sG

pC..: q@s - qC< Y~YzC@ - <bq@S^L zb zPC eCqC^z LCep" z zP- z S - zz- S^ 4Y..SP - ~^S bH b^C%o ^@ HzPC qC..: q@ S <Cq ^b \ b^C%S L- S^C@ bq Ysz>..PSC ebsSfC bq ^CL- zSfC f- YG \ G ^ CszPCq ep" z bq Yss bH b^Cz- q%H^@si yPC \ - fS ~\ C^eCzC@ep" z Hq\ PbY@S^L S^ zP- z S^zCqf- YbHzS CS <- Y~YzC@ - ^@ Hhq zPC \ b@CY zb zq %dC q^ Pb...zb zq @CSz ..SY4C ^C<Gss- q%zb - <PSCfCzPC \ - fS ~\ ep" z ebsS^Y S^ zP- z zS CS^zCqf- Y , @S<b-^z S - eeY@ zb zPC qC..: q@ - <bq@S^L zb zq ^s <zsb^ <bszi

reCS < zC^P^S- YS^@S- zbqf ..CqC <PbsC^ ..PC^ eS^W^L zPC @- z- sCz S^ bq@Cq zb H<Ssz- zC H z-qC C^Q zq <zsb^i Gbq zPS eqf4Y \ @ S%PSLP> Yb... ^@ <Ys^L f- YG ..CqC ~sC@>- Y^L ..SP zPC fbY \ Czq @C@ zPC qY zSfC szq^LzP S^@C^ - ^@ zPC szb^W \ b\ C^z- \ zP- z S \ G s-qC@ fS a 3, i

yPS C^fSip^ \ C^z S Y S@ 4%zPCH<z zP- z S b^%o sG " ^- ^<S YS^@S- zbq zb \ b@CysD ~C^<G>..PSC S^ qC Y" ^- ^<S Yzq @S^L zPC \ - qVz S @CeC^@C^z bHzPC fbY zS%bHzPC..bq@> YW - @SG sC b-z4Cq Wbq - ^ C^zCq^ - Y...: q

### 3.2.2 Transformer Q-Learning

Replay buffer stores transitions  $(s, a, r, s')$  sampled from the environment. The buffer is implemented as a circular queue. The size of the buffer is  $N$ . The buffer is filled with transitions as they are generated by the environment. When the buffer is full, the oldest transition is replaced by the newest one.

The replay buffer is used to sample transitions for training. The transitions are sampled uniformly at random from the buffer. The sampled transitions are used to train the Q-network. The Q-network is updated using the sampled transitions.

#### {i|i|i, yq UCzbcq%R e~z - ^@ p CeY%3~' Cq

The replay buffer is implemented as a circular queue. The size of the buffer is  $N$ . The buffer is filled with transitions as they are generated by the environment. When the buffer is full, the oldest transition is replaced by the newest one.

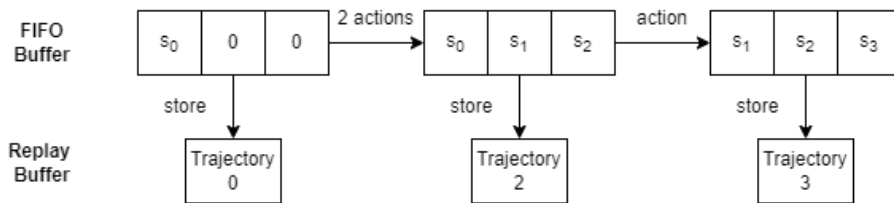
The replay buffer is used to sample transitions for training. The transitions are sampled uniformly at random from the buffer. The sampled transitions are used to train the Q-network. The Q-network is updated using the sampled transitions.

The replay buffer is implemented as a circular queue. The size of the buffer is  $N$ . The buffer is filled with transitions as they are generated by the environment. When the buffer is full, the oldest transition is replaced by the newest one.

The replay buffer is used to sample transitions for training. The transitions are sampled uniformly at random from the buffer. The sampled transitions are used to train the Q-network. The Q-network is updated using the sampled transitions.

The replay buffer is implemented as a circular queue. The size of the buffer is  $N$ . The buffer is filled with transitions as they are generated by the environment. When the buffer is full, the oldest transition is replaced by the newest one.

The replay buffer is used to sample transitions for training. The transitions are sampled uniformly at random from the buffer. The sampled transitions are used to train the Q-network. The Q-network is updated using the sampled transitions.



GS~qC {i|i=? b~4YpCeY%3~' Cq, qPSCz-qC

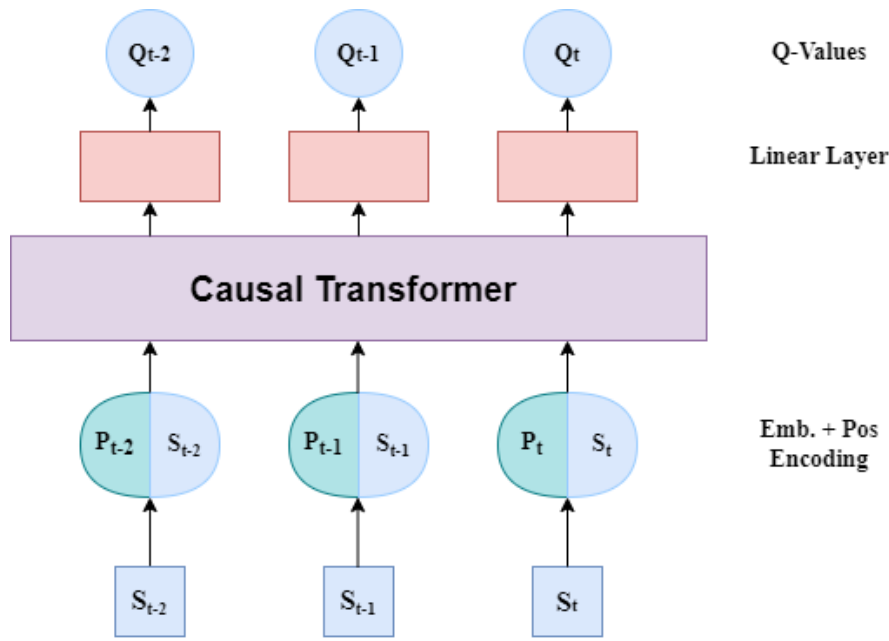
#### {i|i|i3 ] Cz.bqW

The replay buffer is implemented as a circular queue. The size of the buffer is  $N$ . The buffer is filled with transitions as they are generated by the environment. When the buffer is full, the oldest transition is replaced by the newest one.

- ^@zPC bzPCqzPCy- qLz yq ^shbq\ Cq

, 4-zP bHzq UC-zbqCs Ss K@S'zb zPC dbY%yq ^shbq\ Cq. PCqzPCsS'~sb@ YebSSb^- YC\ 4C@S'L ~sC@S' zPC ?y Cf Y- zsb^ Ss ~sC@S' bq@Cq zb qz- S' - sC^sC bHbq@Cq 4Cz. CC^ zS\ CszCes> zP- z - qC ^bz @C^ C@Hbq\ - Y%σ^YW b^ ?yi

yPCs \ C @C-b@Cq 4Y<W' ..SP L-zS'L zP- z ..CqC - eeY@ Hbq S'<qf sC@ sz- 4SS%o. CqC - Yb qCqzSS C@ Hbq zPS CfeCq\ C'z - ^@zPC b-ze-z Y%q qep@<Cs V k Cf YCs zb 4C ~sC@S' zq S'S'L Hb\ - sCz bHV zS\ CszCes zP- z <b\ ebsC- zq UC-zbq%o, s\ - YqCeqCsC^z- zsb^ Ss sPb..^ b^ "L-qC {i|i



$Q_t = \text{Linear Layer}(\text{Causal Transformer}(\text{Emb. + Pos Encoding}(S_{t-2}, S_{t-1}, S_t)))$

$\{i|i|i; yq S'S'L$

yPC qeY%4-' Cq Ss sVb..Y%σ" YC@ - <<bq@S'L zb - ^ QqC@%e bY%σ-s'S'L b^Y%σzPC \ - †S\ ~\ k Cf YC- b4z- S'C@Hb\ zPC Ysz zS\ CszCe bHzPC ebY%σ^Cz. bqWbq Hb\ - q ^@b\ S C@<PbsG

, Hbq b4z- S'S'L C^b-LP zq UC-zbqCs> zq S'S'L sz- qz 4%σb\ e~zS'L zPC \ - †S\ ~\ k Cf YCs Hbq zPC - <sb^Qe- <C bHzPC <-qC^z V sz- zCs fQ<sub>t</sub>g' ..SP zPC ebY%σ^Cz. bqW, Hbq zP- z zPC \ - †S\ ~\ k Cf YCs Hb\ zPC HbYb..S'L sz- zCs - q b4z- S'C@Hb\ zPCz qLz ^Cz. bqW Q<sub>t+1</sub>g

yPCz qLz ..SY4CLSfC^ 4%zPC qC.: q@s b4z- S'C@f<sub>r</sub>t<sub>g</sub>s~\ \ C@zb Q<sub>t+1</sub> >..SP 4CS'L zPC @s<b^z H<z bq yPC Yss H^<Sb^ ..SY4C zPC \ G ^ sl -- qC@ Cqpbq 4Cz. CC^ zPCz qLz - ^@ Q<sub>i</sub>

yPCz qLz ^Cz. bqW..SY<be%zPC e- q \ CzCq bHzPC ebY%σ^Cz. bqW- Hbq - <Cq- S' - \ b-^z bHszCes - ^@zPC CfeYbq zsb^ H<z bq ..SYsYb..Y%σC- %zb CE yPS C^s-qCs sz- 4SS%σ' zPC qC~<qC^z ~e@ zCs zb zPC \ b@Y

yPC HbYb..S'L esC-@bq@Cs~\ \ - qf Cs zPS eqp<Cs=

---

, **Y**l**q**S**P**\ {ic=yq ^s**H**q\ Cqk **Q**C q^S**L**

---

```
Model = CausalTransformer()
targetModel = CausalTransformer()
ReplayBuffer = replayBuffer()
LastStates = laststates()
state = Env:reset()
..PSC timestep < maxTimesteps @b
SHepisodeTimestep < contextLength zPC^
  lastSteps pad(lastStates;state)
CXC
  lastSteps append(lastStates;state)==FIFO
C^@SH
  action = Model:getAction()==E greedyexploration
  s_{t+1};r;d;€ = Env:step(action)
  ReplayBuffer:store(s_t;r;d;a;s_{t+1})
SHReplayBuffer:SamplePossible() zPC^
  QV values = Model(sample(s_t))
  QNext = targetModel(sample(s_{t+1}))
  targets = rewards + (1 - dones) (QNext) )
  Loss = MSE(QV values;targets)
C^@SH
SHtimestep%updateFreq == 0 zPC^
  targetModel:copyParams(Model)
C^@SH
C^@..PSC
```

---

} s^L Lq^Cq q-P S^ bq^Cqzb " ^@zPC 4Csz P%Gq- q \ CzCq <b\ 4S^- zS^ ebssS^Y>zPC **H**Y**B**..S^L P%GqQ  
e- q \ CzCq **H**q zPC zq ^s**H**q\ Cq \ b@CY..CqC <PbsC^=

y- **4**Y { ic= yq ^s**H**q\ Cqk **Q**C q^S^L P%Gq- q \ CzCq

O%Gq- q \ CzCq

, - YC

; b^zCz Y^LzP

I

3- z-P sS C

{ |

3~' Cq sS C

c**CECE**

B\ 4C@@S^L sS C

c|D

Bf- Y- zS^ **H**q ~C^<%o

c**CECE**

**C**E

XG q^S^L q zC

I CQ

OC @s

c

X- %Gq

{

y- qLz ~e@ zC **H**q ~C^<%o

| **CE**

a ezS S Cq

, @ \ „

„ - q\ ~e r zCes

| **CECE**

{ **C**E

### 3.2.3 Transformer Proximal Policy Optimization

R' zPŠ s~4sCzš^ zPC- eeY@ \ CzPb@s zb - @ ez eqđš \ - YebY%bezš S - zš^ S' bq@Cq zb ~sCzq Ū-zbqš - ^@zq ^šHq\ Cq \ b@Cš - qC @SeY%@

yq S'S'L ... s @b^C Hq cEEzPb~s ^@zš CszCes S' - YC^fšp^ \ C^zi yPŠ S' <q sCŠ sš eY%4C- ~sC zPC- \ b~^z bHq S'S'L s \ eYš Š s \ - YCq @-C zb b^QbY% YbqšP \ s - Y. : % ~sS'L HqšP @-z S' szC @bH qC-zšS S'L eqfš~s zq Ū-zbqš - ^@4C- ~sCs- \ eYCC <Š^<%š ..bqCŠ dda i

{i|i{i, yq Ū-zbq%R' e~z - ^@p CeY%š~' Cq

Gbq- ^ - @ ez- zš^ bH ^ b^QbY% YbqšP \ YWC dda >- s \ - Y- \ b~^z bHP- ^Lš - qC ^C-Gs- qš' bq@Cq zb \ - ^z- S' zPCs- \ Cb4Ū-zšfC- s S' zq ^šHq\ Cq k QC q'S'Li

á yPC qeY%4~' Cq ^C@as zb 4C <Y q@- Hq zPC ebY%š ~e@ zC@ S' bq@Cq zb zq S' C<Y-sšfC%b^ ^C. Y%LC^ Cq zC@ zq Ū-zbqš

á yPC YLQđ4- 4SŠŠ eqđ@- <@4% zPC ebY% qC^ b... e- q bHPC qeY%4~' Cq ..PSCb-z<b\ Csz- zš fs\_{t+1}g 4C<b\ C~^ ^C-Gs- qš' yq Ū-zbqš - qC szšYszbq@- s - ..PbYi

á yPCGRGa - qđ%q\ - S's ~^<P- ^LC@- s zPC Ysz V zš CszCes q\ - S' ^C-Gs- qš' bq@Cq zb zq S' zPC ^Cz. bqWzb \ b@CY- sđ ~C^<G

{i|i{i3 ] Cz. bqW

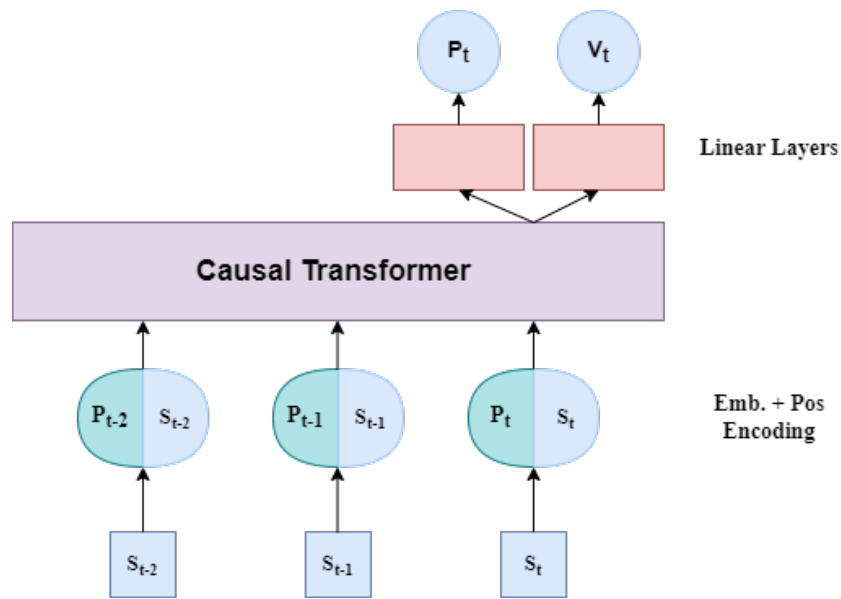
, sP- q@ ^Cz. bqW-sS'L zPC S\ eqđfC@ zq ^šHq\ Cq Hq\ eqfš~s sCzš^ ... s <PbsC^ Hq zPŠ z%@ bH S\ eY\ C^z zš^> qešC^zS'L 4bzP zPC dbY% ^@, - YCH^<š^i BfC^ zPb-LP sC- q zC ^Cz. bqW \ - % - fš@S' zqđq^<C4Cz. CC^ 4bzP bHPCsC b4Ū-zšfC> <Cz- S' <P- q <Cšzšs bHPC^ fšp^ \ C^z \ - %4C bH S' zCšz zb 4bzP sCzš^s bHq S'S'L> \ G ^S'L zP- z - <Cz- S' @LqC bHš \ eYCC <Š^<%š - %4C- <PšfC@ 4% zq S'S'L S' zPŠ HsPš^i

Gbq zPŠ qš sb^ - zq ^šHq\ Cq sš Yq zb zPC b^C ~sC@ S' k QC q'S'L S' S' zq@- <@ ..šP z. b @š Cq^z PC @si yPC f- Y-C PC @ eqđ@- <š - sS'LYC f- YC Hq zPC eqšC^z sz- zC bHPC ^Cz. bqW yPC dbY% PC @ b-ze~zs zPC eqđ4- 4SŠŠ Hq G <P - <š^> ~sS'L - r bH - † - <šf- zš^ H^<š^i yPbsC eqđ4- 4SŠŠ - qC zPC ~sC@ zb <C zC- ; - zLbqš- Y@szq4-zš^ Hq \ ..PCq YLQđ4- 4SŠŠ - ^@C^zpe% qCz- V^ Hq \ i

Gbq zPŠ \ b@CzPCb-ze~zs - qC^bz bHđ ~C^z YHsPš^> sS' <CzP- z bqšL^ - zC@ eqđ4Y\ s ..šP - @f- ^z LC Gzš - zš^ zP- z <b- Y@ ^bz 4C sbYfC@

, 4qšHbqfšC... bHPC ^Cz. bqWš eqšC^z S' " L-q{ij=





$$GSL \sim C \{ i | j = y q \wedge s \text{Hq} \} C q, q \text{PSC} z \sim C \text{Hq} d \text{q} i S \cdot Y d b \text{S} \% a e z S \cdot z b \wedge$$

{ i | i; y q S S L

, S S C @ - \ b \wedge z b H C \wedge f s p \wedge \ C \wedge z C e S b @ C s - c \wedge L C \wedge C q z C @ 4 \% s S L z P C e b \text{S} \% P G @ b H z P C \ b @ C y z b s C C z z P C H \text{Y} b . S L - < z b \wedge i y P C C \wedge z p e \% \wedge @ \text{Y} L Q \text{q} 4 - 4 S \text{S} \text{S} b H e S \text{V} \text{S} L G < P - < z b \wedge - c \wedge s z b q C @ - \text{Y} \wedge L \dots S P z P C b 4 z S C @ c \dots q @ - z C q \ S - Y s L \wedge - Y - \wedge @ z P C Y s z V s z z C s b H z P C \ b @ C y

, H C q C q < P S L z P C \wedge C s s - c \% \wedge \ b \wedge z b H z S \ C s z C e s z P C \ b @ C y . S Y 4 C \sim e @ z C @ R \ b q @ C q z b @ b z P - z z P C q z \sim q \text{q} \text{q} b \text{H} \ G < P z S \ C s z C e \wedge C C a s z b 4 C < b \ e \sim z C @ y P - z S - < P C f C @ S \ z P C s \ C H s P S \wedge b H . P - z \dots s @ b \wedge C S \ z P C ? y S \ e Y \ C \wedge z z b \wedge > 4 - z P C q z P C q z \sim q \text{q} \text{q} b - c \wedge b q \ - \text{Y} C @ - < b q @ S L z b z P C q \ G \wedge - \wedge @ s z \wedge @ q @ @ C f S z b \wedge - < q s s b \wedge C C e S b @ C

y P C \sim c q \wedge z e - q \ C z C p b H z P C \wedge C z . b q \text{W} - c \wedge s z b q C @ - \wedge @ \wedge C \dots \text{Y} L Q \text{q} 4 - 4 S \text{S} \text{S} > C \wedge z p e \% \wedge @ f - Y C s \dots S Y 4 C L C \wedge C q z C @ S \ G < P z q S S L C e b < P > - H C q z P C \ b @ C y S \sim e @ z C @

, @ f \wedge z L G s 4 C z . C C \wedge z P C \wedge C . C q \text{Y} L Q \text{q} 4 - 4 S \text{S} \text{S} - \wedge @ z P C b y @ b \wedge C s \dots S Y 4 C < \text{Y} \sim Y z C @ \sim s S L [ b \wedge z C Q - q \text{Y} C s z S \ - z b \wedge - \wedge @ \sim s C @ z b < \text{Y} \sim Y z C z P C r \sim q \text{p} L - z C X b s s > z P - z \dots S Y 4 C < \text{Y} e C @ 4 C z . C C \wedge [ 1 \quad ; 1 + \ ] z S \ C s z P C q z b s 4 C z . C C \wedge z P C \sim c q \wedge z \text{Y} L Q \text{q} 4 - 4 S \text{S} \text{S} - \wedge @ z P C b y @ b \wedge C s i

, \wedge C \wedge z p e \% 4 b \wedge \sim s \dots S Y 4 C - @ C @ z b z P C Y s s H \wedge < z b \wedge z b C \wedge s - c \wedge C \text{f} e \text{Y} q z b \wedge - \wedge @ z P C f - Y C e b \text{S} \% \dots S Y 4 C < b \ e - c @ \dots S P z P C q z \sim q \text{q} \text{q} b f S \ \ G \wedge q \text{I} \sim \sim c @ C \text{q} \text{p} \text{q} < b \ e b s S L b - q Y s s H \wedge < z b \wedge i

, H C q z P C L q @ C \wedge z @ C s < C \wedge z s z C e > z P C 4 - C q S < Y c @ - \wedge @ z P C e q < C s s s z \text{q} s b f C q \ d s C - @ b q b @ C - \wedge @ s b \ C b H z P C \ b s z S \ e b q z \wedge z \text{H} q \ \sim Y s - c \wedge e q C s C \wedge z C @ 4 C b . =

---

, **YlbqSP** \ {i|=yq ^sHbq\ Cq dda

---

```
Model = CausalTransformer()
oldModel = CausalTransformer()
ReplayBuffer = replayBuffer()
LastStates = laststates()
..PSC timestep < maxTimesteps @b
s_t = env.reset()
lastStates:clear()
lastStates:append(s_t)
Hbq maxEpisodeLength @b
action; logprobs = oldModel.selectAction(lastStates)
s_{t+1}; r; done; = env.step(action)
ReplayBuffer:append(s_{t+1}; reward; done; s_t; logprobs)
Stimestep % updateTimestep == 0 zPC^
Hbq episode @b
RTG = 0
Hbq Reversed rewardBuffer @b
RTG = reward + RTG
C^ @ Hbq
C^ @ Hbq
Hbq ^~\ Beb<Ps @b
logprobs; entropy; values = model.eval(originalSt; originalAc)
ratio = e^{\frac{\logprobs}{originalLogProbs}}
advantages = rtgs - values
surr1 = ratios - advantages
surr2 = clip(ratios; 1 - ; 1 + ) - advantages
loss1 = min(surr1; surr2)
loss2 = 0.5MSE(values; RTG)
loss = loss1 + loss2 + 0.01 Entropy
loss:backward()
C^ @ Hbq
oldModel.copyParameters()
C^ @ SH
C^ @ Hbq
C^ @ ..PSC
```

---

} sS'L Lqf@G q-P S' bq@Cqzb " ^@zPC 4Gsz P@Cp- q \ CzCq<b\ 4S'- zsb^ ebssSfY>zPC HbYb..S'L P%@CqQ  
e- q \ CzCq Hbq zPC zq ^sHbq\ Cq \ b@CY..CqC<PbsC^ Hbq zPs zq ^sHbq\ Cq \ b@CY

y- 4YC {i|= yq ^shq\ Cqdda P%@Cp- q \ CzCp

O%@Cp- q \ CzCq , - YC

---

; b^zCqz Y^LzP	I
B\ 4C@S^L sS C	c D
	<u>CE</u>
XG q^S^L q zC	I CQ
OG @s	J
X- %Gp	{
Ceb<Ps	DE
; Ye p- ^LC	CE
} e@ zC GqQ ~C^<%o	J 1 CeSsb@C Y^LzP
Bf Y- zS^ HqQ ~C^<%o	cCEEE
a ezS\ S Cq	, @ \ „
„ - q\ ~e rzCes	CEEE

---

# 4

## Results

R' zP\$ <P- ezCq zPC qS- Ys bH zC\ ezS L zb S\ eqfCzPC @C<S\$^ zq ^shh\ Cq- ^@zPC b^S' CpX - @ ez-Q  
zb^s bHzq ^shh\ Cq- q<P\$Cz- qS - q' eqS C^ zC@>- lzCq- eqeCq f- S@ zb^ bHzPC bqlS'- Y- q<P\$Cz- qS  
@b^C- ^@sb\ C bH\$z 4bzzC^ C<W - q' S@C^ zS C@  
y PC" ^- YsCz zb^ ..SYs~\ \ - q' C- ^@ @S<-ss zPC b4z- S' C@ qS- Ysi

### 4.1 Original Architecture Validation

R' zP\$ sCz zb^ zPC qS- Ys bHzP\$ S\ eY\ C^z- zb^ - q' sS\ eY%ab\ e- qC@ ..SP zPC 4- sC S\ eY\ C^z- zb^  
S' bq@Cq zb Csz- 4Y\$P - eqeCq 4- sC\$^ C Hh\ zPC C'eGq\ C^zs zb 4C <b^@< zC@ y PqCCz- 4Ys - q' eqfS@C@  
<<bq@S^L zb G <P C^ fSp^ \ C^z zb 4C Cf- Y- zC@ , Y@ z- sCz - q' q' ebqC@ b^ zPCs- \ Cz- 4YI  
R ..SY4C ~sC@ S' zPC HhYb..S^L z- 4Ys zb @C^ bzC zPC \ G ^ ^bq\ - S' C@ qC..: q@ S' - ^ C^ fSp^ \ C^z - ^@  
R \$ ~sC@ zb @C^ bzC zPC sz- ^@ q@ @CfS zb^ 4Cz..CC^ q- ^s S' - ^ C^ fSp^ \ C^zi

y- 4YC Jic= pGs~Ys b^ zPC,, - WcQ? ? b\ - S'

a ~q̄				3-sCR eY\ C^z-zB^			
	[ C@S\ q̄ CeY%o	[ C@S\	[ C@S\ q̄ teCq̄	[ C@S\ q̄ CeY%o	[ C@S\	[ C@S\ q̄ teCq̄	
R	DcH	Di{J	cCEJv	viv	wiCE		cCEc
R	JiDu	ci{I	Cc	{iCE	ciJ		CE

y- 4YC Ji|= pGs~Ys b^ zPCObecq? b\ - S'

a ~q̄				3-sCR eY\ C^z-zB^			
	[ C@S\ q̄ CeY%o	[ C@S\	[ C@S\ q̄ teCq̄	[ C@S\ q̄ CeY%o	[ C@S\	[ C@S\ q̄ teCq̄	
R	DuII	v{	ccciJc	D iu	vuiv		cCEiv
R	JiJu	{i v	ci CE	uiCE	ciCE		ciD

y- 4YC Ji{= pGs~Ys b^ zPCO- YQ PCz P ? b\ - S'

a ~q̄				3-sCR eY\ C^z-zB^			
	[ C@S\ q̄ CeY%o	[ C@S\	[ C@S\ q̄ teCq̄	[ C@S\ q̄ CeY%o	[ C@S\	[ C@S\ q̄ teCq̄	
R	{DiD	J i__	_CE_	{viv	J iv		DviD
R	{	iuD	iIu	CE	CE		ci{

yPC 4-sCR C S eY\ C^z-zB^ ~sC@zb @CfCbe zPS 4b@%bH..bqWS ^bz b^Y%db\ e-q 4YC 4-z bHC s-qe-ssGs zPCbqfLS-Ys eY\ C^z-zB^ bH? G-SS^ yq ^shbq CqP-fsL C-SPCq- PSLPCq\ G ^ cC..:q@bq-Yb..Cqsz-^@q@CfS zB^ 4Cz..CC^ q^si

rS^CCfCq@PSL...:sS eY\ C^zC@S^ zPCs\ CHsPS^ -s? G-SS^ yq ^shbq CqzPC@S Cq^<C4Cz..CC^ \ b@C\ - %b^Y%4C-zq4-zC@zb zPqCC \ -S' @S Cq^<Cs

á ? S Cq^zS eY\ C^z-zB^ bHCqz S H^<zB^s YW\ ~YSQC @-zzC^zB^ bqebSS^ -Y\ 4C@S^Lsi

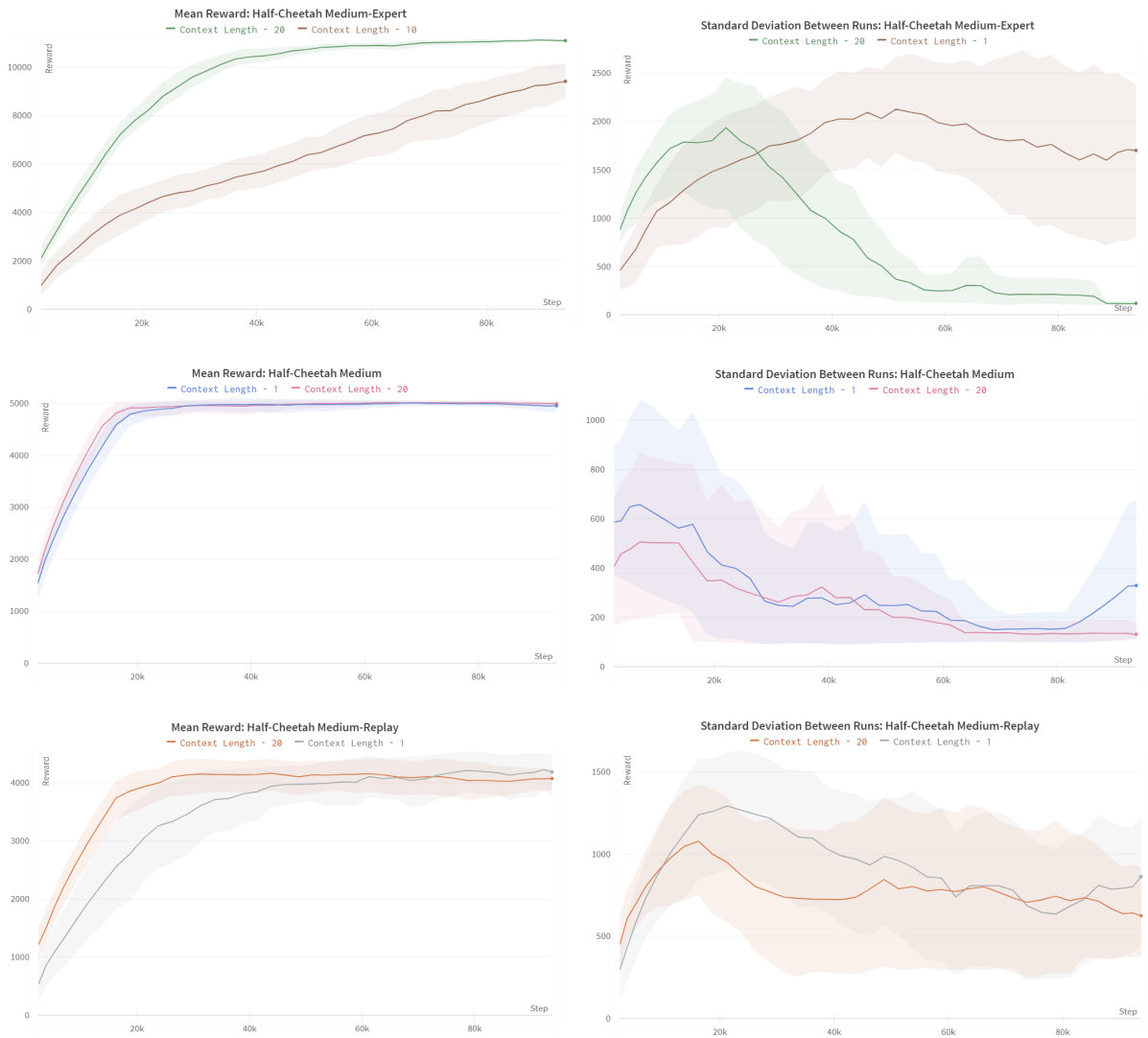
á 3%qCeY<S^L pCX} - <zsf-zB^ H^<zB^s 4%KBX} i

á 3%zq S^S^L zPC\ b@CYf\ s<q z-P S^szG @bHYb @S^L- eqCq S^C@ \ b@CY

#### 4.1.1 Impact of the Context Length

R^ zPS sC-zB^ - ^-zz\ ez-z~^@Cqz-^@S^L..PCzPq|bq^bz S S qG Y%G-Css- q%zb ~sC-b^zCfz bqecfSb-s zS GszGes zb \ b@CY- sD ~C^<CS \ - @G R^ bq@q zb - <PSCfC zPC HYYb..S^L qS~Ys - 4-sCR C \ b@CY...:s zq S^C@..SP - <b^zCfz Y^LzP bH|CE ^@- ^bzPCq ..SP - <b^zCfz Y^LzP bHc>..PSP S zPCs \ C-s P-fsL ^b <b^zCfz-^@sS eY%o.bqWS^L..SP zS GszGes

pGs~Ys - qC@SeY%C@S "L-qC Jici



GSL~cQ Jic= ; b^zCiz XC^LzP= [ G ^ pC..: q@- ^@rz- ^@ q@? CjS zb^ - <qss O- YQ PCz- P ? - z- sCzs

yPC<b^<YsB^s zP- z < ^ 4C@CjC@Hb\ zPCsCqS~Ys sPb...zPCS\ ebqz- ^<CbH-sS'L - <Cqz- S' z/@CbH @ z- sCz ..PC^ zq S'S'L - qS^Hq-C^ z YC q^S'L \ b@CY4%σ-sS'L s-eCqjS@ YC q^S'Li ? SzS'L~sPS'L G <P @ z- sCz - Yb..s HqzPC^Hb..S'L U@LC^ C^zs=

á [ C@S\ QteCq= Rz 4C<b\ Gs G sCq zb \ b@CY<b\ eYs- zC@ zq UC<zbqS ..SP - 4SLCq <b^zCizi yPC \ bq<b\ eYci zPC C^fSp^ \ C^z zPC \ bq S\ ebqz- ^z Sz 4C<b\ Gs zb e-%σ zC^Sb^ zb eqCfSQ b-s zS CszCesi

á [ C@S\ = ? SeY%σ HszCq <b^fCqL^C- ^@s\ - YCq sz- ^@ q@CjS zb^ ..SP <b^zCizi yPC zq UC<Q zbqS - qC ^bz ebYSP@ C^b-LP HqzPS @S Cq^<C zb 4C<b\ C sL^S <- ^z S' zPC C^ qS~Ys ..PC^ <b\ e-q@..SP [ C@S\ QteCq @ z- sCzsi

á [ C@S \ QCeY%oy PCzq UCzbqCs zb qCep@-C- qCsb sS eYzP- z zPC \ b@CY<- ^ - <PSCfCzPCs- \ C  
 <b^<YsS^s ..SPb-z <b^zCfz>4bzP S \ C ^ qC.: q@- ^@S sz ^@ q@CfS zS^ - <qss q^si

yPC@- z- sCz l ~- Y%@C b^szq zCs zP- z zPCS\ ebqz ^<CbHb^zCfz S^<C sCs ..SP zPC <b\ eYfS%bHZPC  
 sbY-zS^i } ^@Cqz- ^@S^L - sD ~C^<CS \ ~P \ bCfS ebqz ^z zP- ^ ~^@Cqz- ^@S^L - zS\ CzCe - ^@zP- z S  
 zPC \ bsz S ebqz ^z- seCz bHC Y^L ..SP @S <-Y C^fSp^ \ C^zi

### 4.1.2 Target Modeling

R^ zPs s~4sC-zS^>- fS-- Y- zS^ bHPb...- zq S^C@ \ b@CYq <zs zb - qD ~CzC@z- qLz qCs~Y S @SeY%  
 S bq@Cqzb Cz- 4YsP zPC 4CP- fS-qbHZPC- LC^z>..PC^ H<S^L @S Cq^z b4C<SfCs S^ zPCs- \ CC^fSp^ \ C^zi  
 yPs @C b^szq zCs SZPC- LC^z S <- e- 4Y bH- ^@Cqz- ^@S^L Ss C^fSp^ \ C^z b-zs@C bH- @Cq\ S^C@  
 Lb- Y- ^@ ..PCzPCq bq ^bz bfCq^ zS^L b<<-qS S^ zPs Hb\ bHC q^S^Li yPC \ -S^ eq\ S<CbHZPS C^eCfS C^z  
 S zb l ~CzS^ ..PCzPCq bq ^bz S - ^ - LC^z qC Y%sbY^S^L - ^ C^fSp^ \ C^z fS @bCs ^bz ~^@Cqz- ^@Pb...zb  
 - <PSCfC \ ~YsY b4C<SfCs S^s@C bHSi

GSL~qC Ji | @SeY%@Pb... @- z- sCz bHC <P l ~- Y%eCqHq s - <<bq@S^L zb zPCs- \ Cz- qLzsi



GSL~qC Ji | = y- qLz [ b@CS^L - <qss O- YQ PCz- P ? - z- sCz

R' zPC \ C@S \ Q†eCq @ z-sCz - ^ S'zCqSs'L ePC^b \ C^b^ P-eeC^si ? ~Czb Ss Y<WbH4- @l ~ Y%o  
 q^> zPC \ b@CYb^Y%o~^@Cqz- ^@s e- qz bHzPC LSfC^ z-dCzi yPS S CfsC^z - <bq@S'L zb zPC seSVC S  
 l ~ Y%zP-z P-eeC^s ..PC^ zPCz-dCz - <PSCfCs - 4b-z 50% bHSs \ - †S \ - \ f-Y-G  
 yPS ..b-Y@4C@C^ ^C@-s- ^ bfCq^ zS'L 4b~^@ q%oS^<CzPC \ b@CY@bGs ^bz~^@Cqz- ^@..P-z S P-eeC^  
 4C%b^@zP-z z-dCz - s @SeY%G@S " L~cC Ji{i



GSL~cC Ji{i = a fCq^ zS'L 4b~^@ q%oS^ | C@S \ Q†eCq @ z-sCz

rS^<CzPC \ C@S \ @ z-sCz Ss - ^ CfC^ \ bCq qC@-<C@ @ z-sCz zP- ^ zPC \ C@S \ Q†eCq @ z-sCz zPC  
 Ys Ss bHSs ~^@Cqz- ^@S'L - cC CfC^ \ bCq qC@-<C@> sPb..S'L zP-z zPC \ b@CYP- s - Y bsz ^b qG <zb^ zb -  
 <P- ^LCs^ z-dCzi yPS ..b-Y@ G ^ zP-z zPC pyK zq S^S'L sSL^ - Ys CfC^ ~^^C-Gss- q%oS^ zPS \ b@CY>sS^<C  
 zPCqC Ss ^b @Ss^<zb^ 4Cz..CC^ Lbb@- ^@4 @ q^s S^sCzPCs \ C@ z-sCz

R' zPC \ C@S \ Q†eY%@ z-sCz zPC sS~zS^ <P- ^LGs sS^<CzPC \ b@Cys \ eYs - qCeY%4~' Cq zP-z  
 H<C@zPC eqLqSsSb^ bHLbS'L Hb \ <Cq ~^@Cqz- ^@S'L bHzPC C^fSp^ \ C^z zb - H Y~^@Cqz- ^@S'Li Gbq  
 zP-z qG sb^ - ^@zP- ^W zb zPCeqS^<CbHSfCps%zPC \ b@CYP- s - ^ - Y bsz Y^G q4P- fSb-q- <fss z-dCz  
 b4U<zSfCs>4CS'L bHb-qC Ys S@4%δbz Wb..S'L - ^%C†eCq sqz zLCS Hbq zPC 4Csz sbY-zb^i

Rz Ss G s%zb <b^<Y@C zP-z zPC \ C@S \ Q†eY%@ z-sCz LSfCs zPC \ b@CY- 4CzCq ~^@Cqz- ^@S'L bHSs  
 C^fSp^ \ C^z zP- ^ - ^%bzPCq@ z-sCz - ^@zP-z Hbq zP-z qG sb^ zPCS^zq@-zb^ bH bC@SfCps%zb PSLPCq  
 l ~ Y%@ z-sCz Ss - ...%zb sbYfC zPS @SfCps%4bzzY^C^Wb~^@S' PSLPCq l ~ Y%@ z-sCzi

## 4.2 Architectural Changes

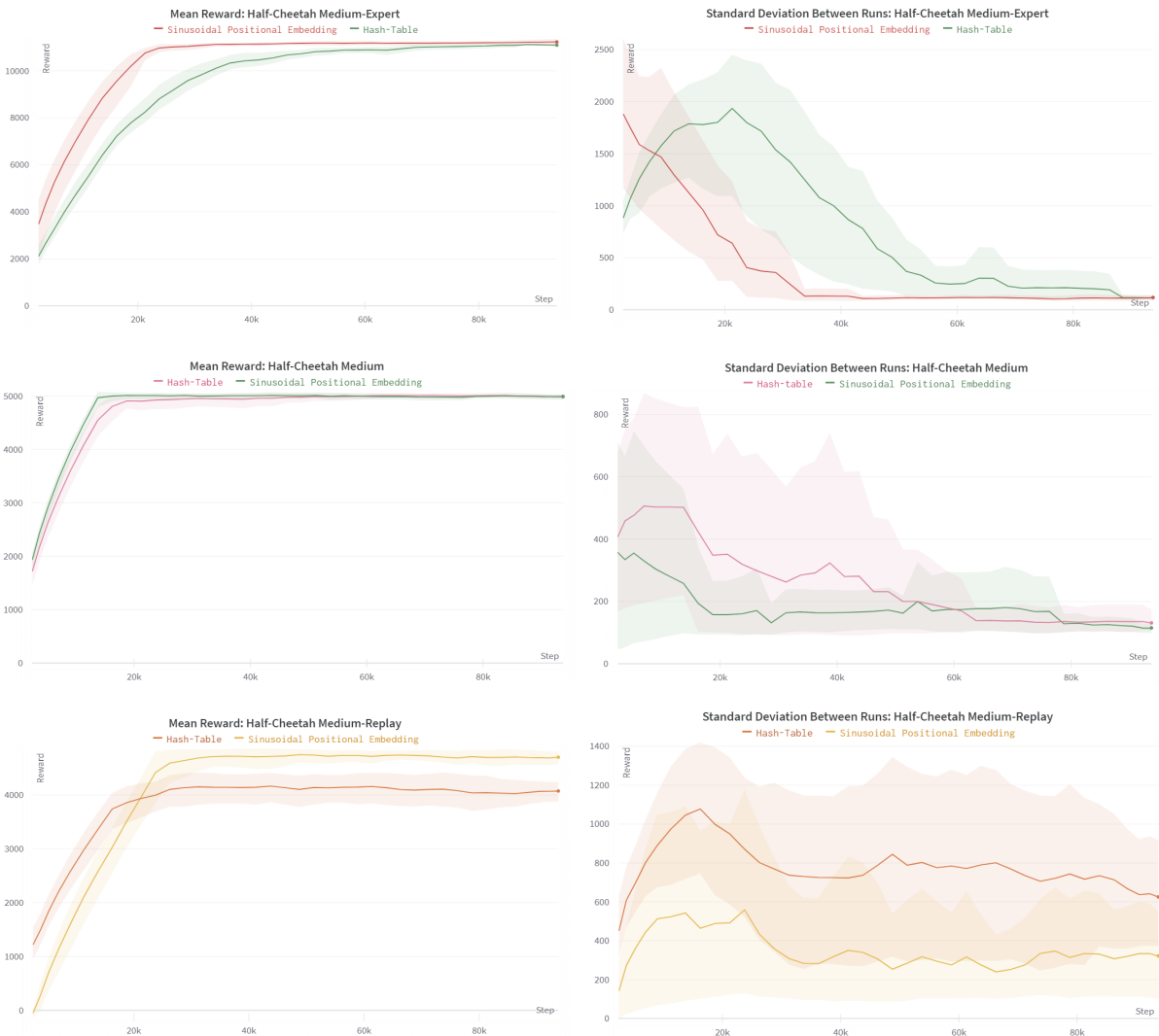
R' zPS sCzb^ zPC qS~Ys bHeCqHq \ S'L \ ~YsY<P- ^LGs zb ? y - qCeCqS^zC@

yPC \ Czfs ~sC@ zb - ^- Y%C \ b@Cz z VC Szb - <b~^z zPC \ G ^ qC..: q@ b4z- S^C@ - ^@ Ss sz- ^@ q@  
 @CfS zb^ 4Cz..CC^ q^si rS^<C? y Ss - \ CzPb@ zP-z ~sCs s-eCqfS^C@ YC q^S'L> zPC zq S^S'L Yss Ss - Yb  
 z VC Szb - <b~^z>CfC^ zPb-LP Sz ..SY^bz 4CzPC \ - S^ \ Czfs bHCf Y- zb^i



## 4.2.1 Positional Embeddings

, s eqfSb-s%a\ C^zB^C@ebsSb^- YC\ 4C@S^Ls ..Gq P- ^LC@Hb\ P-sPQ- 4Ys zb sS~sbS@ YebSbSb^- Y  
 C\ 4C@S^Ls zP-z -qC YC q^- 4YC - ^@ ^bz Hb\ C^i yPC HbYb..S^L qS~Ys ..SY 4C eqS^zC@ - <bq@S^L zb  
 C^fSp^ C^z - ^@@-z-sCz yPCsq - ^- %sS ..SY 4CLYb4- YS^ bq@Cq zb S@C^zH%4CP- fSb-q Ye-zzCq^si  
 GSl~qC Jij qeqS^zS zPCb4z- S^C@ qS~Ys S^ zPCO- YQ PCCz- P @b\ - S^i



GSl~qC Jij= dbsSb^- YB\ 4C@S^L= [ G ^ pC..q@- ^@r z ^@ q@? Cfs zb^ - <qss O- YQ PCCz- P ? - z sCz

Rz Ss G s%zb <b^<Y@CzP-z - <bq@S^L zb zPC S^<q S^L <b\ eYq^S%bHzPC @ z- sCz ~sC@zPC\ bqCS Q  
 ebqz ^z Sz 4C\b\ Gs zb ~sC- \ bqC-b\ eYq^ z%CbHebsSb^- YC\ 4C@S^Lsi XC q^- 4YCS^~sbS@ YebSbSb^- Y  
 C\ 4C@S^Ls - Yb...Hbq •C^SYC \ bC@S^L bHLS GszGes - ^@- 4CzCq @S^z<zb^ 4Cz.CC^ zPC 4CLS^S^L - ^@  
 C^@bHYb^L CeSb@Csi

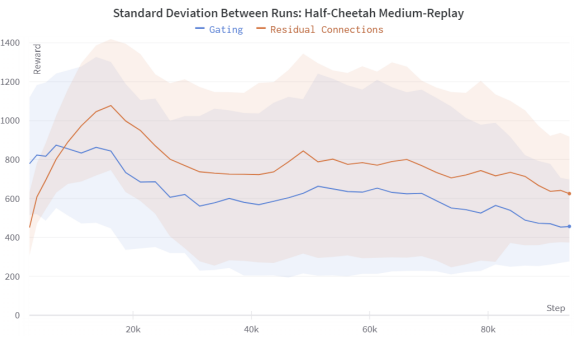
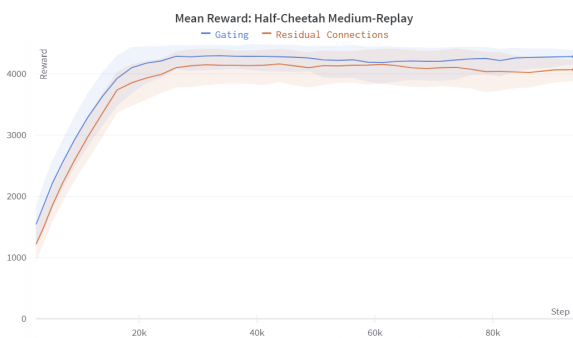
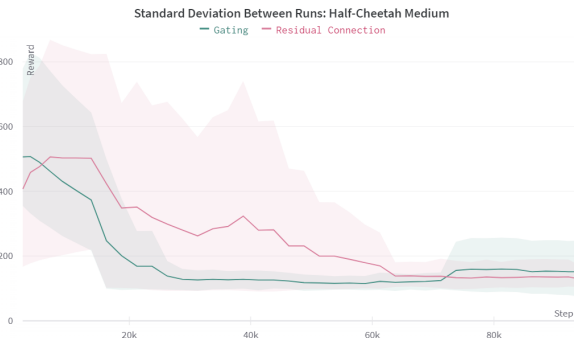
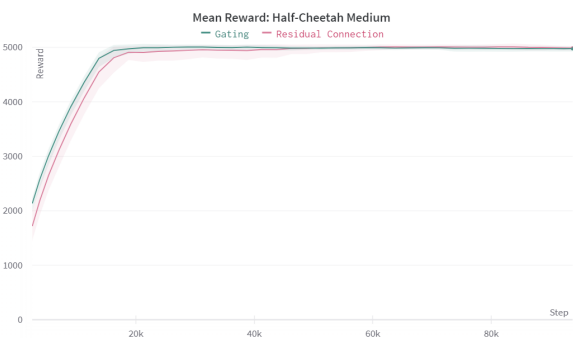
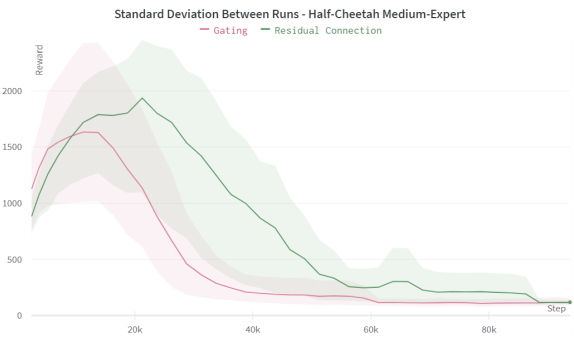
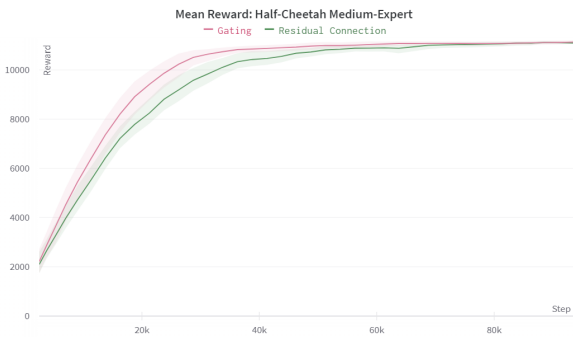
B- q%ob^ fCqL^ <G> Yb. Cq sz ^@ q@ @CfS zsb^ - ^@ - ^ bfCq YS^ <qG sC S^ sz- 4S%o qC sb\ C bHzPC S\ ebqz ^z zq Ss Hb^ @- <qbss C^ fSb^ \ C^zsb^ - ^ S eY\ C^z- zsb^ bHC q^ - 4YsS^ ~sb@ YebSSb^ - YB\ Q 4C@S^ Lsi yPS @bCs ^bz qC^ Cz b^ - sL^S < ^z S eqfC\ C^z b^ zPC C^ @qS~ Ys bHzPC\ b@CY 4-z Sz @bCs PCé zPC\ b@CY 4G b\ C\ bCfS- \ eYCC <S^z- ^@sz- 4Y.. PC^ ~sS^L \ C@S- \ bq\ C@S- \ Qf eCqz @z- sCzi BfC^ zPb-LP zP- z <b^ fCqL^ <C seCC@ S P-qz b^ zPC\ C@S- \ QCeY% @z- sCz> zPC bfCq Y eCqHq\ - ^ <C 4C^ C zs LqG zY% d\ ~sS^L sS^ ~sb@ YebSSb^ - YC\ 4C@S^ Lsi yPC qG sb^ Hq zPS \ - %AC qY zC@ zb zPC H-z zP- z - YC q^ - 4Y ebsSSb^ - YC\ 4C@S^ L \ - %dC q^ e- zCq's zP- z - qC ^bz bzPCq. SC eqS C^z S^ zPC @z- sCz> . PS-P - Yb..s Hq- s\ - Y- \ b^z bHCi zq ebY zsb^ d\ - ^ bzPCq. SCsS eY @z- sCz zP- z eq@- <Gs \ C@b <qC qS~ Ys .. PC^ <b\ e- q@ .. SP zPC q\ - S^S^L bezb^ si

### 4.2.2 Gating

„ SP zPC b4U <zSfC bH^ ^@S^L - sz- 4Y sbY zsb^ > - zC\ ezs zb S^SS YC zPC \ b@CYS^ - ... %zP- z .. b-Y@ 4CP- fCsS\ SY q%zb - [ - qMf? C~Sb^ d q <Css \ G ^z- <P- ^LCS^ zPC\ b@CY- qPSzCz- qCzP- z .. b-Y@ ^b... S^ <Y @CL- zS^L \ <P- ^S\ s S^ szC @ bHqS@ - Y b^ ^ C zsb^ si

y PCHbYb.. S^L qS~ Ys > zP- z q\ - S^ <bPCq^z - <qbss @z- sCzs qfC YzPCS\ e- <z zP- z L- zS^L zq- Y%P- s b^ zq S^S^L > zP- z - Yb zq ^sY zC zPC\ sCyfCs b^ - ^ S^ <qG sC bHe- q \ CzCq 4% zPC zq ^sHq\ Cq ^Cz. bqWbH11% > LCzS^L \ bsz bHzPC - qPSzCz- qS S^ zPC | \ SYb^ e- q \ CzCq q ^LG

? S Cq^ <Gs 4Cz. CC^ 4bzP - qPSzCz- qS - <qbss @z- sCzs S^ zPC O- YQ PCCz- P @b\ - S^ - qC qCeCqS^ zC@ S^ " L~qC Jili



$GSL \sim qC Jil = K - zS^L = [ G \wedge pC : : q@ - \wedge @ rz \wedge @ q@ ? CfS zB^ - < qss O - YQ PCCz P ? - z sCs$

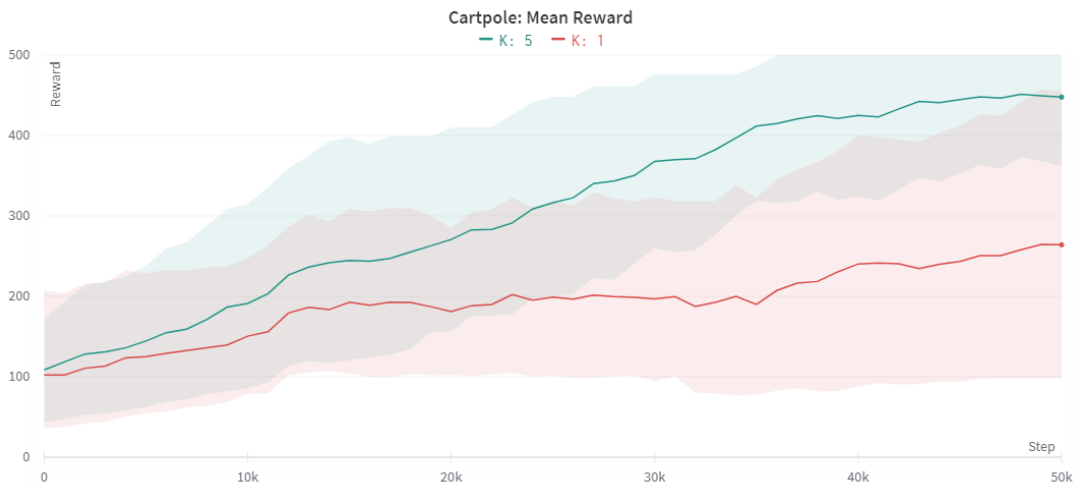
BfC^ zPb-LP zq S^S^L...b-Y@z W- zbYS^ zC ebq YC <S^<%@-CzbzPCS^<q sCS^ e-q \ Cq>P- fS^L S^<q sC@ 4Cz .CC^ 35 - \wedge @ 42% Hbq - s^LYC q^> zPC eb .GfHYeGfHbq - \wedge < C S\ eqfC C^zs S^ sz \wedge @ q@ @CfS zB^ - \wedge @ sz-q@S^C ss zb q \wedge @b\ S^SS Y^ - zB^ < - \wedge \wedge bz 4C \wedge @Cqz- zC@

; b^fCqL^<CS s^LPz%HzCq - \wedge @CfC^ SHz ...s^bz zPCq eS@ @CqG sCS^ sz \wedge @ q@ @CfS zB^ ..b-Y@- z Y^sz \ G^ zP- z zPC \ b@CY^b...4CP- fCs \ bC<bPCq^zY%<qss q^s \ G^S^L zP- z zPCeG WeCfHbq - \wedge <Cs \ - %4C<b\ Cs^LPz%b .Cq 4-z zP- z zPC ..bqz q^s 4C<b\ CY^ss bH^ Ss-G

rz \wedge @ q@ @CfS zB^ ..b-Y@- \ b-\wedge z zb P-\wedge @C@S bHqC : : q@ ebS^zs S^ sb\ Cq^s - \wedge @ zPS S^ \wedge b Y^LCq S^ zPC < sC ..SP L-z^L> \ G^S^L zP- z zPC S^SS Y^ - zB^ bHzPC^Cz .bqW..SP L-z^L \ <P- \wedge S^ \ s qC Y% zq^s Y^zCs S^ - 4CP- fS^>q sS^ Y^q zb - [ - qWf ? C^SS^ d q<Cssi

### 4.3 Transformer Q-Learning

, Hq - eeYSL zPC - @ ezC@ k XC q^S'L - YbqP\ - ^@ " ^Cq~^S'L zPC ^C<Gss- q%P%Gpe- q \ CzCp> zPC  
 \ b@CY...s Cf- Y- zC@S' zPqCC^ fSp^ \ C^zs - ^@<b\ e- qC@..SP Sz s <b~^zCq- qz ..SP ^b <b^zCqz S' bq@Cq  
 zb eqpeCq%~^@Cqsz- ^@HzPCeqG^<CbH sD ~C^<C\$ qCCf- ^z zb zPCsbY-zb^ bHzPS z%CbHeq4Y\ si  
 R' zPS <- sCb^Y%qC... q@..SY4C@SeY%G@sS^<Csz- ^@ q@<CfS zb^ S @Gz-zY%sb^GzC@zb zPC<b^zCqz  
 Y^LzP S' zPCs \ C...%s zPC \ G ^ qC... q@ yPC \ -S' ebS'z zb qCz- S' S zP-z zPC<b^<Cez bHb' QbYS%  
 Y q^S'L ..SP sD ~C^<C\ b@CS'L <- ^ 4C- eeY@- ^@sbYfCeq4Y\ s S' pX...SP LqG z C <C^<%s sPb...^  
 S' " L~qC Jiv

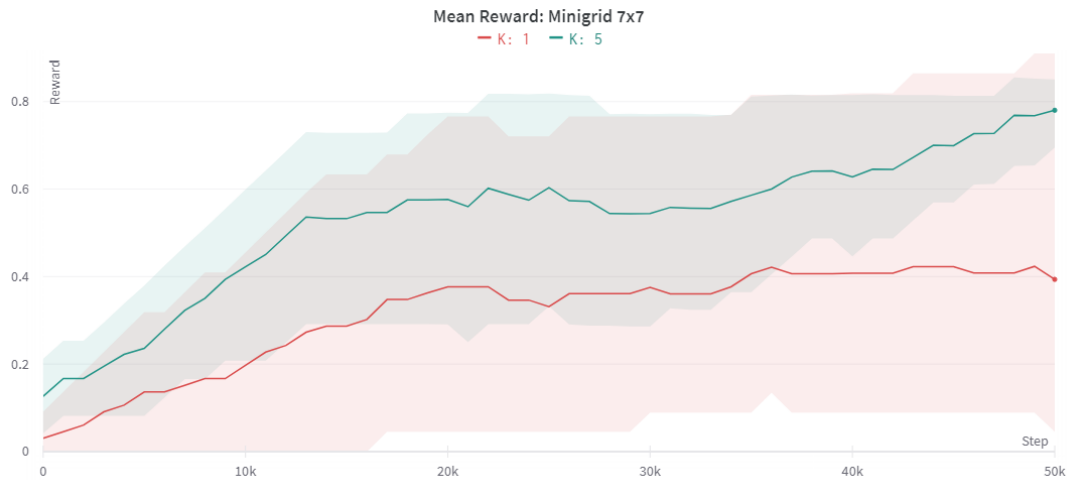


GSL~qC Jiv= ; - qebYpGs~Ys Hbqyq ^shbq Cqk QC q^S'L

R' zPC ; - qzdbY C^fSp^ \ C^z ^bz b^Y%SS - sbY-zb^ -<PSCfC@ S' cEzPb-s- ^@ zq S'S'L szCes> 4~z -  
 \ - fS ~\ s<bqC S G sY%b4z S' C@ ..SPS^ - s\ - YzC\ ebq Y4b~^@ q%..PSC qCz- S'S'L sz-q@S' Gss S' G <P  
 Cf- Y- zb^i yPS eqfCs zP-z ^bz b^Y%SS Sz ebss4Y\ zb ~sCyq ^shbq Cq b^ b' QbYS%Y q^S'L> 4~z zP-z  
 zPC%o qC- Ysb- ^ C <C^z \ CzPb@zb YC q'i

„ SP - ^ S<qC sC@- \ b~^z bHzq S'S'L szCes- ^ [ Xd bq- zq ^shbq Cq..SPb-z <b^zCqz...b~Y@C' ^S%  
 <b^fCqC zb zPCs \ C qS~Y> 4~z zP-z ..b~Y@ eqfC ~^H sS4Y\ S' qC Yeq4Y\ s ..PqC YC q^S'L b^ zPC Lb  
 - ^@ Hsz S Gss^zS Y

R' [ S'Lq@>zPC~sCbHe- qCqC... q@s eqp@-<G sS\ SY qCqS~Ys> eqfS'L zP-z zPC \ b@CY<- ^ 4C- L^bszS  
 zb Ss qC... q@H^<zb^>- s zPC qS~Ys S' zPC [ S' SK q@ 7x7 C^fSp^ \ C^z @SeY%S " L~qC Jivi

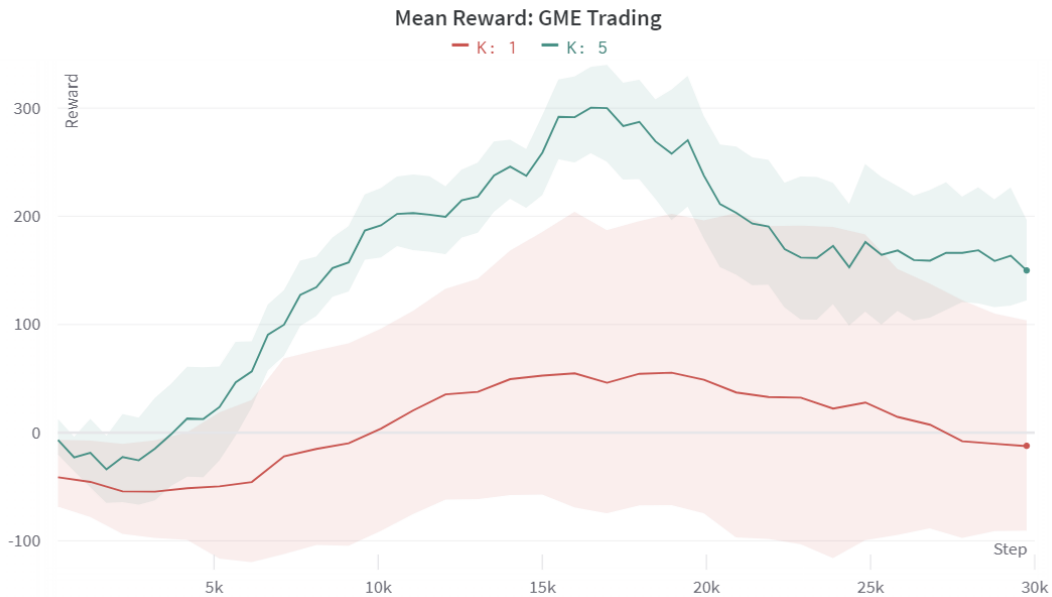


GSL~qC Jiu= [ S^Lq@pGs~Ys Hhqq ^shhQ Cqk QC q^S^L

yPCL-e S^ YG q^S^L 4Cz.CC^ - \ b@CY..SZP - ^@..SZPb~z <b^zCiz S\$ \ bqCfS@C^z S^ zPS\$ sCzS^L>..PGqC  
 - \ b@CY..SZP <b^zCiz S\$ - 4Yzb " ^@zPCsbY-zSb^ zb zPC\ - <C80% bHzPCzS\ Cb^ - fCq LG>..PSC- \ b@CY  
 ..SZPb~z <b^zCiz S\$ b^%o 4Yzb @b S 40% bHzPCzS\ G yPS\$ S\$ ..P-z zPC\ G ^ qC..: q@\ G ^s S^ zPS\$ <sC  
 sS^<C- ^bq\ - YS @qC..: q@ bHc S\$ eqS^C^zC@ CfCq%zS\ CzPC C^ fSb^ \ C^z S\$ sbYfC@

GS^ Y%oS^ zPC zq @S^L C^ fSb^ \ C^z zPC \ b@CY 4LS^s 4%o YsS^L \ b^C%dq eGfHq\ S^L S^ - ^G-zq Y  
 \ - ^^GqLq @- Y%YC q^S^L zPC <fS- YebS^s bHb^L - ^@sPbq - <S^s S^ - sD ~C^<G B fC^z-- Y%zPC  
 <b^zCizQ...: qC\ b@CYC q^s zb b4z- S^ - eG Weq" z bHc |\_iDj zS\ Gs zPC bqfLS^ - YS^ fGz\ C^z b-zeGfHq\ S^L  
 zPC f- YCbHs eY%PbY@S^L zPCszb <W- ^zS^s eG Wf- YG>..PS^P ...: s { CEs\ Gs Szs bqfLS^ - Yf- YG>..PSC zPC  
 \ b@CY..SZPb~z <b^zCiz b^%b4z- S^C@ - eq" z bHccIi|c zS\ Gs zPC bqfLS^ - Yf- YG

yPC YG q^S^L qS~Ys - q@SeY%@S " L-qJiDi



Mean Reward: GME Trading

Figure 4.4: Mean Reward: GME Trading. The graph shows the performance of two different configurations (K=1 and K=5) over 30,000 steps. The K=5 configuration consistently achieves a higher mean reward, peaking at approximately 300 around step 16,000, while the K=1 configuration peaks at around 50 around step 18,000. Shaded areas represent the variance of the mean reward.

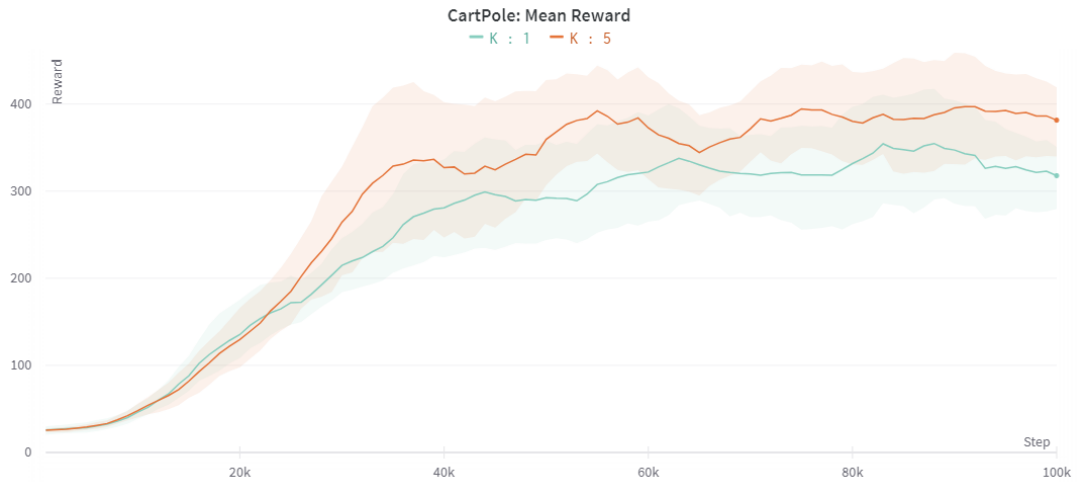
The K=5 configuration shows a significant advantage in terms of reward, likely due to its ability to capture more complex patterns in the data. The K=1 configuration, while simpler, struggles to maintain a positive reward over time, eventually dropping below zero.

#### 4.4 Transformer Proximal Policy Optimization

The Transformer Proximal Policy Optimization (PPO) algorithm is used to train the agent. It involves a series of updates to the policy and value functions, with a focus on minimizing the variance of the gradients. The algorithm is implemented using the following code:

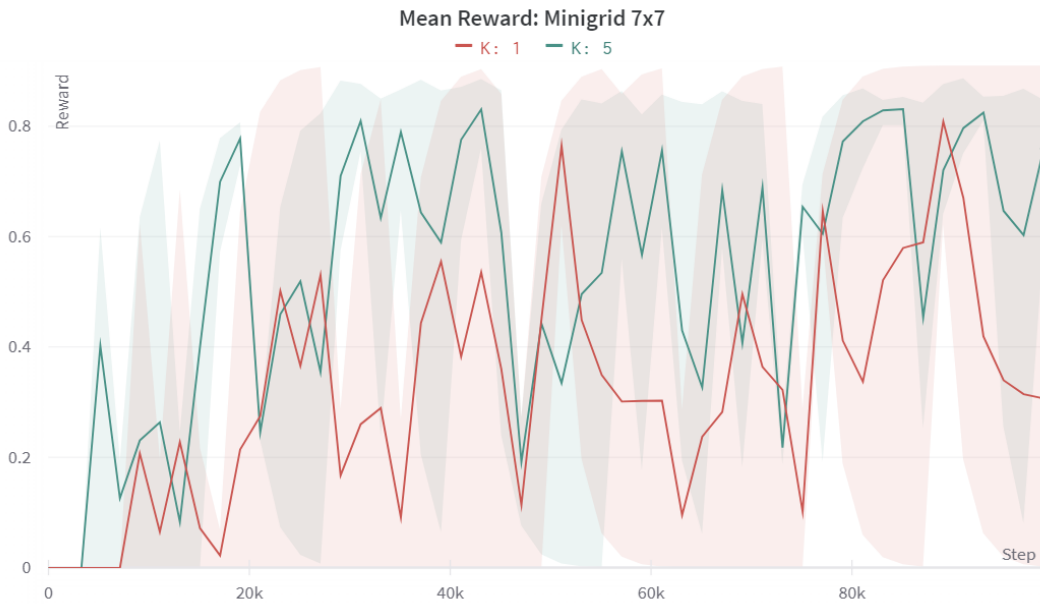
```
Observe env.reset()
for episode in range(10000):
    obs, _ = env.reset()
    done = False
    total_reward = 0
    while not done:
        action = policy(obs)
        obs, reward, done, _ = env.step(action)
        total_reward += reward
    value_loss = loss(policy(obs), total_reward)
    value_loss.backward()
    policy_optimizer.step()
    value_optimizer.step()
    policy_optimizer.zero_grad()
    value_optimizer.zero_grad()
    total_reward = 0
```

The code above shows the training loop for the Transformer PPO algorithm. It involves a series of updates to the policy and value functions, with a focus on minimizing the variance of the gradients. The algorithm is implemented using the following code:



GSL~qC Ji\_ = ; - qebXCpGs~Ys Hhyyq ^shhQ Cqdda

„ PC^ H<S'L - se- qPC qC.: q@H^<Sb^>zPCb^QbS%o YbqzP\ szq-LLYs zb YC q'i KCzS'L - qC.: q@ 4C<b\ Gs - q qC sS-- zS^>sS^<CqG <PS'L zPC Lb- Y- <bq@S'L zb q ^@b\ - <Sb^s Ss @S <~Yi yPSs 4C<b\ Gs - ee- q^z S^ " L~q JicE...PCqC YC q^S'L Ss ebsSAYC JfC^ S^H - qW@4%o^sz- 4C qCs~Ysi

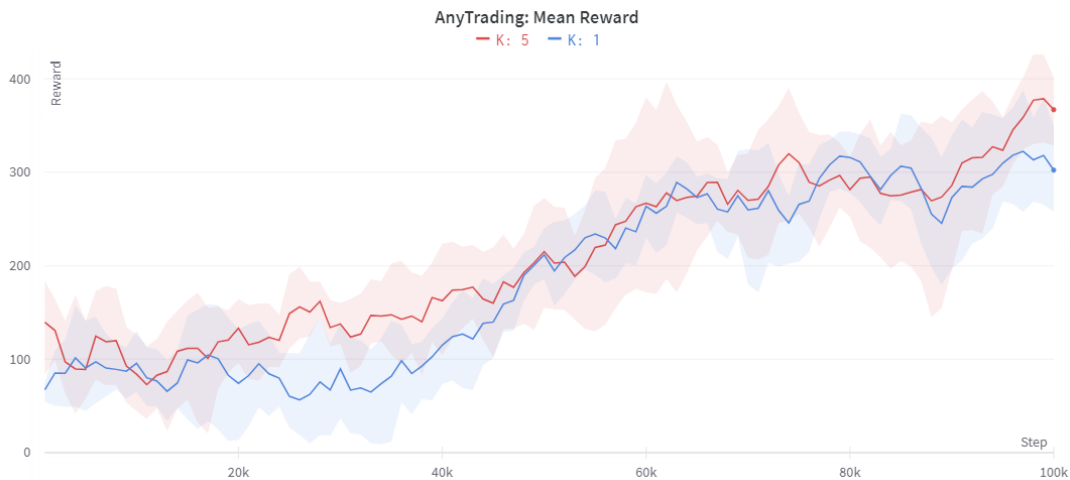


GSL~qC JicE [ S^S q@pGs~Ys Hhyyq ^shhQ Cqdda

a ^<CzPCLb- YS qC <PC@>zPCzq UCzbcq%zP- z...s z- W^ <- ^bz 4Cs- \ eY@fS qeY%4-' CqCqCg zC@%o \ G^S'L zP- z zPC\ b@YzC^@s zb P- fC@S <~YSS LCzS'L C^b-LP YC q^S'L sS^<- Ys zb HY%oLq se zPC Lb- Y yPSs \ - Ws zPC\ b@YLCz sz~<W~^Yss zPC ~^YWC%ePC^b\ C^b^ bH^ ^@S'L zPC Lb- Y\ ~YSe YZs Gs 4%o <P- ^<CP- eeC^si

3C- ~sCzPS \ b@CY-sS- \ bCqCfeqGssSfCqC... q@H^<zb^ zP- ^ zPCeqfSb-s C^fSip^ \ C^zs>zPCqS- Ys  
 b^ zPC- ^%mq @S'L C^fSip^ \ C^z 4C<b\ C\ bCqCeq\ SS'L zP- ^ b^ k QC q^S'L>\ - VSL Sz- Lbb@4C^<P\ - qW  
 zb Cf- Y- zCzPCyq ^shhq\ Cqdda - @ ez- zS^i

GSL~qC Jicc sPb..s zP- z zPC\ b@Cyb4z- S's- sYLPz%PSLPCqC... q@..PC^ ~sS'L <b^zCfz zP- zS^ bz fCq%  
 \ G^S'LHY..PSP Ss Ss S'Yqzb zPC 4CP- fSbq b4sCqfC@S' ; - qzdbYI Ob..CfCq 4C- ~sCzPCqC... q@H^<zb^  
 S^bz S'G q...SP eq^ z>zPS \ CzPb@%SCY@C@- \ - †S \ - eq^ z bHcl C|c zS Gs zPC bqfLS- Yf- YC...PSC  
 ^b <b^zCfz b^%P- @- eq^ z bHcl i|c zS Gs zPC bqfLS- Yf- YG



GSL~qC Jicc=, ^%mq @S'L pGs~Ys Hbqyq ^shhq\ Cqdda

## 4.5 Discussion

yPSs sC-zS^ eqsC^zC@zPC qS~Ys bHcuc C^fSip^ \ C^z zq S^S'L - ^@ Cf- Y- zS^ sCssb^s - <qss v C^fSip^Q  
 \ C^zs - ^@{ seG-S < - YbqzP\ s zP- z ..Gc S eY\ C^zC@- ^@ zq S^C@ Hb\ s<q zP S^ bC@Cq zb eqeCq%  
 f- Y@ zCzPC<-qq^z sz- zCbHzPC- q- ^@zb Cfe- ^@~eb^ Si

, ^ C bCz zb Cf- Y- zC CfCq%GSL^ <PbSC- ^@zb UszS%CfCq%CfeCq\ C^z 4CHC Ss qC YS- zS^ ...s  
 @b^C S^ bC@Cq zb b4C%zb s<C^zS < qfLbq a zPCq CfeCq\ C^zs zP- z @S^ ^bz <b^fC%ebsSfC qS~Ys ..Gc  
 - Yb @CfCbeC@YWCzPC~sCbH b^, Sy 9\_>r..S^ yq ^shhq\ Cq 9 CE>p Cy zSfCdbSzb^ - YB^<b@S Ls 9{: - ^@  
 CfC K, | s 9{:i

Rz < ^ 4C <b^sCqC@zP- ^ 96.5% bHzPCfeCq\ C^zs ..GcS~<<GssHYsS^<CzPC Cf- Y- zS^ b^ dda ..SP  
 zq ^shhq\ Cq S^ - se- qCSczS'L eq@<C@Hqzbb ~^sz- 4C qS~Ys zb eqeCq%eb^<Y@C- ^%PS Li

, - Y@ zS'L zPC? GSS^ yq ^shhq\ Cq eqfC@zP- z Ss qS~Ys - qC qeCq@<SfC- ^@zP- z zPCqC- qC szSY  
 sb\ C 4bzXC^G W zb z<WCS b^ S^C qS^HqC\ C^z YC q^S'L>- ^@zP- z zPCqC Ss szSYqbb\ Hbq S eqfC\ C^z  
 ..PC^ @GSL^S'L zPC ^G-q Y^Cz..bqW zP- z C\ eY%zPC\ CzPb@s bHzPC bqfLS- Ye- eCq

yq ^shhq\ Cq k QC q^S'L - ^@ dCf†S - YdbS%a ezS S- zS^ - @ ez- zS^s ..Gc s~<<GssHY- s eqbHq



<b^<Cez-eeY< zB^s bHebe~Yqb' QbY%o^@b^QbY%oYlqzP\ s S' sD ~C^<C\ b@CS'L...SP YqCpb\ HqS\ eqf\ C^z- ^@ebzC^zS YS' zPC- qG i

XbbVSL -z z 4Y JiJ>- <b\ e- qsb^ 4Cz.CC^ 4bzP - YlqzP\ s sPb...zP-z zPC%o qC <b\ eCzSfCS' - Y C^fSp^ \ C^zs - ^@zP-z sbY-zb^s zb - Yeq4Y\ s - qC - Y.: %Hb-^@>CfC^ SHCss zP- ^ S@C Y

y- 4Y JiJ= pG~Ys; b\ e- qsb^ bHa ^S^CpCS HqC\ C^z XG q^S'L...SP rD ~C^<C\ b@CS'L

	dda	? k ]
; - qzdbY	JDI i{	I CE
[ S^Lq@fsbyC%g	DE	DE
, ^%mq @S'L fep" zg	c _iD]	cI CE  c

# 5

## Conclusions and Future Work

„ Szp zPCLb YbH-sS'L zq ^shbq\ Cq zb sbYfC qS' HbqC\ C'z YC q'S'L eq4Y\ s .. SzPb-z zPC Y\ Sz-zb^ bH  
” †C@@ z- sCz> zPS @ssCqz-zb^ - ^- Y%C@\ CzPb@s zP-z z-<WCzPCS\ e- <z bHkCqz- S' - q-PSzCz-q Y-Pb&Cs S'  
sz- 4SS%o ^@s- \ eYC C <SC^<%Jfb\ .. PCqsb\ C \ CzPb@s eqjfC@ CseC-S Y%o-sCHYfsS'-sb@ Yebsszb^ - Y  
C\ 4C@@S' Ls - ^@L- zS'L \ G-P- ^S\ sgi yPS sz-@%o Yb C\ eP- sS C@b^ zPCS\ ebqz- ^<C bH@SfCqS%oS' zPC  
<qG-zb^ bH@ z- sCzs bqzPC^CC@ Hbq Cfe Ybq zb^ zb <qG-zCs- S@ @SfCqS%o

? G-SSb^ yq ^shbq\ Cq ...s ~sC@-s - sz- qzS'L ebS'z zb - <PSCfC zPC eqpebsC@ b4UC-zSfCs> eqjS'S'L -  
sz- qzS'L ebS'z zb ~^@Cqz- ^@ Pb... sD ~C^<C \ b@CS'L - ^@b" S'C qS' HbqC\ C'z YC q'S'L .. CqC S'zCqzb^ Q  
^G-zC@ Rss b4UC-zSfCs .. CqCzq ^sYzC@ zb b^ YC qS' HbqC\ C'z YC q'S'L - ^@zPC- YbqzP\ S- - <PSCfC\ C'zs  
-qC- LqG z eqbH@b^ <Cez bH..P- z <^ 4C- <PSCfC@ S' zPS- qG i

### 5.1 Achievements

„ Szp zPC @CfCb\ C'z bHzPS @ssCqz-zb^> zPC HbYb.. S'L b4UC-zSfCs .. CqC- <PSCfC@=

á R@C^zS' <-zb^ bHsd ~C^<C \ b@CS'L 4bzzXC^Cw S' zPC? G-SSb^ yq ^shbq\ Cq - q-PSzCz-qC- ^@zPCSj  
@SfCz <b^^G-zb^ zb zPC ~sC@ @ z- sCzi ; YsC@ @ z- sCzs - qC @C\ C@ eq4Y\ - zS- - ^@ @SfCqS%o ^@

CfeYbq zsb^ - qcepebs@- s- sbYzsb^ zP- z- Yb..s Hbqsz-q%XC q'S'L zP- zS ^bz sz~<Wb^ HY Yb^L - sS'LYCb4UCzsfG yPCqCS - Y.:%- 4- Y^<C4Cz.CC^ bezs\ - YeCqHbq\ - ^<CS - ^ b4UCzsfC- ^@LqC- z eCqHbq\ - ^<CS \ ~YSeY>4-z- ^ C'fSip^ \ C'zS \ - @CbH \ ~YSeYsYLPz%@S Cq'z b4UCzsfG>..PSP CfeYS's zPCS\ ebqz- ^<CbHzPS zbeSi

á Bf- Y- zsb^ bH@S Cq'z zq ^sHbq\ Cq- q-PSzCz-q Y<PbsCs S' b" S'C q'S'Hbq\C\ C'z YC q'S'Li K- zS'L \ C-P- ^S\ s Lq ^z - @@@sz- 4S%P- ^W zb zPCqS'SS YS- zsb^ - s - [ - qWfS ^ ebY%..PSC YC q'Q - 4YC sS~sb@- YebssSb^ - YC\ 4C@S'LS LqC z%sd eqfjC ~^@Cqsz- ^@S'L bHY^L sD ~C^<Cs - ^@- Yb... Hbq\ bqC • C'fSIP%P- ^ - " †C@z- 4Y

á , eeYs- zsb^ bHzq ^sHbq\ Cq b' \ b@CqC b' QbY% ^@ b^QbY%qS'Hbq\C\ C'z YC q'S'L - YbQ qfP\ s>\ - W^L \ CzPb@bYLSs Hbqsd ~C^<C\ b@CS'L S' pX- f- SY 4YC- s - sz- q'S'L ebs'z zb @CfCbe HqPCq..bqWS^ zPC- qG i

## 5.2 Future Work

yPS 4b@%bH.bqW..b-Y@LqC z%4C^C' z Hb\ @CfCbe\ C'z S' zPC HbYb..S'L @qCzsb^s=

á Rz Ss ebsS4YC zP- z CfeYbq zsb^ - ^@LC^Cq YS- zsb^ \ - %P- fC- ^ S'zCqb^ ^Czsb^ ..SP sD ~C^<C\ b@C CS'Li R eqfjC\ C'zs S' zPS ebs'z \ - %PCe <qC zC- LC^zs zP- z- qC <- e- 4YC bH \ ~YSe- sW^L - ^@ HszCq- @ ez- zsb^i

á dqCq S'Co\ b@CS \ - %o Yb...Hbq- s\ - YCq- \ b~^z bHzq S'S'L - ^@sS eY%Hz-q @CSL^ <PbsCsi a" S'C eqCq S'S'L S'zb b' QbY%qz S'S'L \ - %o <PCfC- ^ C <S^z <b\ 4S'- zsb^ ..PC^ z- <W^L s- \ eYCC <S^<%o

á [ bqC <b\ eY† C'fSip^ \ C'zs ^CC@zb 4C~sC@zb Cf Y- zCzPCeqpebs@- YbqfP\ si ? ~Czb <b\ e~Q z-zsb^ - Yb S- zsb^s C'fSip^ \ C'zs YWCzPC, z- qf 4C^P\ - qWbq qG Yq4bzS C'fSip^ \ C'zs zP- z ~sC - <z- - zbs ..CqC ^bz Cf Y- zC@

á yPCyq ^sHbq\ CqdqfS\ - YdqfS\ - Ya ezS S- zsb^ ^CC@s zbeq@~<CV f- YGs - ^@V YbLQq4- 4S%CS bHzPC- <sb^ se- <Czb 4C4 <VQqpe- L- zC@ a ~q\ CzPb@bYL%4Cb\ Cs YS SC@zP- ^W zb b^%oS'L zPCe- sz- ^@eqsC'z zbeq@szb^C- <sb^ - ^@f- YC- HqCzPC<-qq^zsz- zC>..PC^ zPCe- sz eq@szsb^s <b- Y@ 4CqC ^C@

á GS- Y%o.bqY@ \ b@CS>sD ~C^zS YfCqsb^s bH\ ~<Cq 9J:>B' <S^zCq 9I: - ^@bzPCq \ b@CS zP- z ..CqC ^bz ~e@ zC@zb sD ~C^zS Y b@CS \ - %4C CfeYbq@zb H Y%o^@Cqsz- ^@ ..PCzPCq bq ^bz zPCqC Ss - Yb qb\ HbqS eqfjC\ C'zs S' zPbsCsCzS'Lsi

# Bibliography

- ☉: Xi ; PC^>Vi X->, i p- Uš..:q ^>Vi XCC>, i KqfCq [ i X-sVŠ>di , 44CCY, i r qŠ Sf-s- ^@R [ bqQ @z<P> ? C-SŠb^ zq ^shhql Cq p CS Hhqc C^z YC q^S^L fS sD ~C^<C\ b@CS^L> CE | CEci
- ☉: Oi f-^ O-ssCZ^ i ? bcp^>Gi r zq-4> [ i OGssCZ ] i rb^>Cq z- ^@Ti [ b@%SY ? Cce qCS Hhqc C^z YC q^S^L - ^@zPC@C @%azS @ | CEI a ^ŠCi , f-SY 4Y= Pzse=ww qf:Sfibqlw 4swcDc|iCEvJD
- ☉: pi ri r~zbb^ - ^@, i Ki 3-qzb^>pCS Hhqc C^z XC q^S^L= , ^ R^zpb@<zb^> | ^@Cai [ Ry eqSs> | CED>S 3] =uDCEDu(CECEci
- ☉: yi K~ez>, i V-\ -zP>, i VC\ 4P-fS- ^@?i ObsC > yb..:q@s LC^Cq Ye~qbsC fSŠb^ s%szC s> | CEci
- ☉: di OC^@Cqb^> pi RY\ > di 3-<P\ -^> Ti dS^G~> ?i dqC~e> - ^@ ?i [ CLCq ? Cce qCS Hhqc C^z YC q^S^L zP-z \ -zzCq> ; bpp> fbY - 4swcuCEiCEIvCE | CEui a ^ŠCi , f-SY 4Y= Pzse=ww qf:Sfibqlw 4swcuCEiCEIvCE
- ☉: ^ i y-ss->^ i ? bcp^>, i [ ~Y@ Yyi BqC>^ i XS ? i ; -s-s>?i 3~@C^>, i , 4@bX - YVŠ Ti [ CqCZ , i XCq^ ^< i y i XSŠ q e- ^@ [ i pSC@ SYCq ? Cce\ S@ <b^zpbYs-SZ> CE | CEI
- ☉: Xi V-Sq [ i 3-4-CS-@CP> di [ Sbs> 3i a sŠ sVŠ pi Oi ; -\ e4CY Vi ; <CPh..sVŠ ?i BqP- ^> ; i GS^> di Vb- V..sVŠ ri XCfS^C pi rCe-ssŠ Ki y~<Vq - ^@ Oi [ SP-YC.sVŠ [ b@CQ 4-sC@ qCS Hhqc C^z YC q^S^L Hh q -z qŠ ; bpp> fbY - 4swc\_CeiCEIu> | CE\_i a ^ŠCi , f-SY 4Y= Pzse=ww qf:Sfibqlw 4swc\_CeiCEIu
- ☉ yi 3i 3q..^> 3i [ - ^> ] i p%CCq [ i r~44SP> Ti V-eY^> di ? P-qŠ.:Y , i ] CCYW^z ^> di rP%o\ >Ki r-szq%o, i , sWY ri , L-q.:Y , i OCq4CqQ bss>Ki Vq-CLCq yi OC^SLP- ^>pi ; PS@ , i p-\ Gp>?i [ i ŠSLYCq Ti „ ~>; i „ S^zCq; i OGssC [ i ; PC^>Bi rSLYCq [ i XSz..S^>ri Kq %o 3i ; PGss> Ti ; YqW ; i 3Gq^Cq ri [ < - ^@ŠP> , i p-@Hq@ R r~zVqfCq - ^@ ?i , \ b@CS X-^L~LC \ b@CS - qC IC. QPbz YC q^Cq> ; bpp> fbY - 4sw| CEIcJcvI> | CEI a ^ŠCi , f-SY 4Y= Pzse=ww qf:Sfibqlw 4sw| CEIcJcvI

9: Ti ? Cj^S>[ i ; P- ^L>Vi XCC>- ^@Vi yb-z- ^bf> 3Bpy=eqCq S^S^L bH@Cce 4S^S^C-zS^ -Yzq ^sQ  
Hbq\ Cq Hbq Y^L-- LC~^@Cqz- ^@S^L> ; bpp>fbY -4swcDcCEEDCE>|CED

9E, i , -s.- ^S | i rP- <Cq | i d- q- q- Ti } s^WqCz> Xi Tb^Cs> , i | i Kb\ G> Xi V- S^Cq- ^@  
R dbYs~VPS^> , zzC^zS^ S -Y%~ ^CC@ ; bpp>fbY -4swcuCviCEuv|>|CEui Q ^Y^Ci , f- S^4Y= Pzes=ww qj: SfibqLw 4swcuCviCEuv|

9c: pi Gi dq-@C^<S> [ i pi ai , i [ -†S b> - ^@ Bi Xi ; bY\ 4S^S , s-qfC%ob^ b^ S^C  
qS^Hq-C C^z Yq^S^L= y- †b^b\ %o qfSC.> - ^@ beC^ eq4Y\ s> |CEi Q ^Y^Ci , f- S^4Y= Pzes=ww qj: SfibqLw 4sw| CEiCE{ Du

9|: , i rPCqzS^sW%o G- ^@\ C^z- Y bHqC~qC^z ^C-q Y^Cz.bqWfp] | g- ^@ Y^L sPbqCq\ \ C bq%o  
fXry[ g ^Cz.bqW dP%SS- ?=] b^Y^G q dPC^b\ C^> fbY JCE> ei c|{CE> \ -q |CECE Q ^Y^Ci , f- S^4Y= Pzes=ww qj: SfibqLw CfcCEvh | GueP%@| CE\_ic|{ CE

9{ : Vi KqC > pi Vi rCf- sz- f> Ti Vb-z^ W 3i pi rzC- ^C4q^W - ^@ Ti r<P\ S^P-4Cq  
Xry[ = , sC-qP se- <C b@SSC%o ; bpp> fbY -4swcI CEiCECE\_> |CEi Q ^Y^Ci , f- S^4Y= Pzes=ww qj: SfibqLw 4swcI CEiCECE\_

9J: Vi ; Pb> 3i f- ^ [ CqPC^4bCq- èi K YCPqC> Gi 3b-L- qC> Oi r<P..C^W- ^@ ^ i 3C^Lb> XG q^S^L  
ePq sC qCeC^z- zS^s ~s^L p] | C^<b@Cq@C<b@Cq Hbqsz- zS^s- Y\ - <P^Czq ^sYzS^> ; bpp>fbY -4swcJCEicCED>|CEJi

9I: , i ? bsbjSzsW%oXi 3C^Cq , i VbYC^SWf>? i , C^SSC^4bq^†i ŠP- S-yi } ^zCqP^S^Cq [ i ? CPLP- ^S  
[ i | S^@CqCqKi OCLbY@ri KCY%Ti } s^WqCz>- ^@] i Ob~Y4%o , ^ S -LC S ..bqP cv†cv ..bq@s=  
yq ^shbq\ Cq Hbq S -LC qC-bL^S^S^ -z s<- YC ; bpp>fbY -4sw| CECEcc\_|>|CECE Q ^Y^Ci , f- S^4Y= Pzes=ww qj: SfibqLw 4sw| CECEcc\_|\_

9v: , i , q^-4> [ i ? CPLP- ^S Ki OCLbY@ ; i r~^> [ i X- <S> - ^@ ; i r<P\ S@> , S^S= ,  
fS@b fS^S^ zq ^shbq\ Cq ; bpp>fbY -4sw| cCEicIv\_c> |CEci Q ^Y^Ci , f- S^4Y= Pzes=  
ww qj: SfibqLw 4sw| cCEicIv\_c

9u: Ti , Y\ \ -q yPCSY-szq zC@ zq ^shbq\ Cq |CED Q ^Y^Ci , f- S^4Y= Pzes=ww Y^Y\ \ -qLS^P~4iS^w  
SY-szq zC@Cq ^shbq\ Cq

9D Vi OC†i ŠP- ^L>ri pC^>- ^@Ti r~^> ? Cce qS@- YC q^S^L Hbq S -LC qC-bL^S^S^> ; bpp>fbY -4swc|c|CE{D|>|CEi Q ^Y^Ci , f- S^4Y= Pzes=ww qj: SfibqLw 4swc|c|CE{D|

9\_: ri @, s<bY^Oiyb~fcp^>[ i Xi XG fSzz> , i ri [ bqbs>Ki 3SpY^S- ^@Xi r-L~^> ; b^fS= R eqfS^L  
fS^S^ zq ^shbq\ Cq ..S^P sbH <b^fbYzS^ -YS^@- <S^C 4S sC> ; bpp>fbY -4sw| cCEicCE\_u> |CEci Q ^Y^Ci , f- S^4Y= Pzes=ww qj: SfibqLw 4sw| cCEicCE\_u

9CE Ši XS> ^ i XS> ^ i ; -b> Oi O-> ^ i „ CS Ši ŠP-^L> ri XS> -^@ 3i K~b> r..S zq ^sHbq\ Cq  
OSq q-PS- YfSS^ zq ^sHbq\ Cq~sS'L sPŠC@..S@b..s ; bpp> fbY -4sw|cEicJCE|CEci @ ^ŠCi  
, f-SY4Y=Pzes=ww qf:SfibqLw 4sw|cEicJCE

9c: , i p-@Hq@ Ti „ ~> pi ; PS@ ? i X-^> ? i , \ b@CS R r~zVfCq Cz -YŠ X^L~LC \ b@CS -qC  
~^s-eCqfSSC@ \ ~YŠz sWGC q^Cq> a eC^, R4bL> fbY c>^bi D>ei \_>|CE\_i

9|: di 3Yb\ > yq ^sHbq\ CqŠ Hb\ s<q zP> CE |CE\_i

9{ : Ši ?-S Ši ^-^L> ^ i ^-^L> Ti Ki ; -q^b^CY> ki , i XG-^@ pi r-YVP-z@Šbf> yq ^sHbq\ CqQY  
, zC^zSfC Y^L~LC \ b@CS 4C@^@ - " †C@C^LzP <b^zC†z> ; bpp> fbY -4swc\_CeiCEiDvCE |CE\_i  
@ ^ŠCi , f-SY4Y=Pze=ww qf:SfibqLw 4swc\_CeiCEiDvCE

9J: , i [ ^SP> Vi V-f-W-bLY> ? i rSfCq , i KqfCs R , ^zb^bLYb-> ? i „ S@szq> -^@  
[ i pS@ SYCq dY%SL -z qŠ ..SP @Ce qS^HqC C^z YC q^S'L> |CE{i @ ^ŠCi , f-SY4Y=  
Pzes=ww qf:SfibqLw 4swc{c|iivCE

9I: yi rS b^S^S ? Ce lQC q^S'L ..SP se-<C S^f-@Cq> |CEci @ ^ŠCi , f-SY4Y= Pzes=  
wwP~LLS^LH<G<bw4YbLw@CeQYQ^ ^

9v: Oi f^ O-ssCY>, i K-G>-^@? i rSfCq ? Ce qS^HqC C^z YC q^S'L..SP @b-4YCl QC q^S'L> |CEi

9u [ i Ti O~sWGCpz-^@di rzb^C> ? Ce q<-qC^z lQC q^S'L Hbq e-qS Y%b4sCqf-4Y \ @es> ; bpp>  
fbY -4swc|CEiCEi|u>|CEi @ ^ŠCi , f-SY4Y=Pze=ww qf:SfibqLw 4swc|CEiCEi|u

9D Ti r<P-Y-^>Gi „ bYV>di ? P-qŠ.-YŠ , i p-@Hq@-^@ai VŠ bf> dq†S -YebY%bezS S-zš^  
-YbqP\ s> |CEui @ ^ŠCi , f-SY4Y=Pzes=ww qf:SfibqLw 4swc|CEiCEi{Ju

9\_ : Ti r<P-Y-^>ri XCfS^C>di [ bqz>[ i R Tbq@^>-^@di , 44CY> yq-sz qLS^ ebY%bezS S-zš^>  
|CEi

9CE , i [ ^SP> , i di 3-@S> [ i [ S†-> , i KqfCs yi di XSYq e> yi O-qC% ? i rSfCq  
-^@ Vi V-f-W-bLY> , s%<Pq^b-s \ CzP@S Hbq @Ce qS^HqC C^z YC q^S'L> ; bpp> fbY  
-4swc|CEiCEiD>|CEvi @ ^ŠCi , f-SY4Y=Pze=ww qf:SfibqLw 4swc|CEiCEiD

9c: Vi ; b44G> Ti OSz^> ai VŠ bf> -^@ Ti r<P-Y-^> dP-sŠ ebY%Lq @S^z> ; bpp> fbY  
-4sw|CEiCEJcv>|CECE

9|: Bi d-qšbzb> Oi Gi rb^L> Ti „ i p-G>pi d-s<^~>èi K YCPq>ri [ i T-%W\ -q [ i T-@Cq4Cq>  
pi Xi V-~H-^>, i ; YqV ri ] b-q%[ i [ i 3bzfS^S^W ] i OCCs>-^@ pi O-@sCY> rz-4SS^S'L  
zq ^sHbq\ CqŠ Hbq qS^HqC C^z YC q^S'L> ; bpp> fbY -4swc\_cCEuvJ> |CE\_i @ ^ŠCi , f-SY4Y=  
Pzes=ww qf:SfibqLw 4swc\_cCEuvJ

9{ : Oi Gi rb^L>, i, 4@bX - YVb>Ti yiref^LC^4Cq>, i; YqW Oirb%GpTi,, ip-Gri] b-φ% i, P-U>  
ri XS->? iy S<- \ -Y>] i OCCss>? i 3Cb f>[ i, ip S<@ SYCp- ^@[ i [ i 3bz fS^SW , Q da=b^Q bY%  
\ - †S \ - \ - ebsz Cfb cfebY%bezS S- zS^ Hq @S< qzC- ^@ <b^zS^-b-s <b^zcpY S DP R^zCq^-zb^- Y  
; b^HqC^<C b^ XC q^S^L pCe qS^C^z zS^s> R Xp | C E , @S , 4 4 > Bz PbeS >, eqSY|vQ C E | C E  
a e^p C f S . i ^ C z > | C E E @ ^ Y S Ci , f- S Y 4 Y C = Pzes=wwb e^ C q f S . i ^ C z w Hq \ n s E r % a Y J G f O

9J: Ti r <P \ S P - 4 C p p C S Hq C C^z Y C q^S^L ~es S C @b. ^ = ? b^ z eq C @ s z q C . : q @ s Q U - s z \ - e z P C \ z b  
- < z S ^ s > ; b p p > f b Y - 4 s w c \_ c | i C E D u l > | C E \_ i @ ^ Y S Ci , f- S Y 4 Y C = P z e = w w q f S f i b q L w 4 s w c \_ c | i C E D u l

9I: ? i r S y C p , i O ~ - ^ L > ; i [ - @ @ s b ^ > , i K ~ G > X i r S q > K i ? q C s s < P C > T i r < P q f z . S S C p R , ^ z b ^ b L Y >  
, i d - ^ ^ C q s P C y f - \ > [ i X - ^ < z b > r i ? S C X - ^ > ? i K q C . C T i | P - \ > ] i V - Y P 4 q C ^ C p R r - z s W f C p  
y i X S Y S q e > [ i X G < P > V i V - f - W - b L Y > y i K q G e C Y - ^ @ ? i O - s s - 4 S > [ - s z C f S L z P C L \ C b H L b  
.. S P @ C e ^ C - q Y ^ C z . b q W - ^ @ z C C s G q p > / - z - q > f b Y I | \_ > e e i J D J J D \_ > C E | C E v i

9v: [ i T - ^ ^ C p k i X S - ^ @ r i X C f S C > a " S C q S Hq C C^z Y C q^S^L - s b^C 4L s C ~ C ^ C \ b @ C S L  
e p 4 Y C > S , @ f ^ < S S ^ / C - q Y R ^ Hq \ - z b ^ d q < S s S ^ L r % z C s > | C E c i

9u: [ i p C @ ^ i ^ - \ - @ > - ^ @ r i K ~ > ; - ^ .. S W C @ S P C e b " S C q S Hq C C^z Y C q^S^L m C E | C E i

9D , i 3 - ^ S b > , i d i 3 - @ S > T i ; i ,, - W C p y i r < P b Y C > T i [ S c p f S > - ^ @ ; i 3 Y ^ @ C Y ; b 4 C q Y  
; b^zq szSfC 3Bpy Hq qS^HqC C^z Y C q^S^L > ; b p p > f b Y - 4 s w c C E i C E J { c > | C E c i @ ^ Y S Ci  
, f- S Y 4 Y C = P z e s = w w q f S f i b q L w 4 s w c C E i C E J { c

9\_ : k i Š P C ^ L > , i Š P - ^ L > - ^ @ , i K c p f C p a ^ Y S C @ C S S b ^ z q ^ s Hq \ C p S d q < C C @ S L s  
b H z P C { \_ z P R ^ z C q ^ - z b ^ - Y ; b ^ H q C ^ < C b ^ / - < P S ^ C X G q ^ S ^ L > s C q d q < C C @ S L s b H [ - < P S ^ C  
X G q ^ S ^ L p C S C q p > V i ; P - ~ @ P - q S r i T C L C W > X i r b ^ L > ; i r . C e S f q S K i | S > - ^ @  
r i r - 4 z b > B @ s i > f b Y c v | i d [ X p > c u | { T - Y | C E | > e e i | u C E | | u C E \_ i @ ^ Y S Ci , f- S Y 4 Y C =  
P z e s = w w e p < C C @ S L s i \ Y i e c C s s w f c v | w P C ^ L | | < P z \ Y

9CE Bi y b @ b p f > y i B q C > - ^ @ ^ i y - s s > [ ~ U b = , e P % S s C ^ L S ^ C Hq \ b @ C q - s C @ < b ^ z c p Y S / C E /  
R B B w p r T R ^ z C q ^ - z b ^ - Y ; b ^ H q C ^ < C b ^ R ^ z C Y L C ^ z p b 4 z s - ^ @ r % z C s > | C E | > e e i I C E v I C E { i

9c: Ti G - > , i V - \ - q a i ] - < P - \ > K i y ~ < W C p - ^ @ r i X C f S C > ? J q Y ? - z - s C z Hq @ C e @ z - @ q f C  
q S Hq C C^z Y C q^S^L > | C E E

9| : p i , L - q . : Y ? i r < P ~ ~ q \ - ^ s > - ^ @ [ i ] b c p ~ < S , ^ bezS S z S e C p e G z S f C b ^ b " S C q S Hq C C^z  
Y C q^S^L > | C E \_ i @ ^ Y S Ci , f- S Y 4 Y C = P z e s = w w q f S f i b q L w 4 s w c \_ C E i C E J {

9{ : R Ti K b b @ H Y b . > Ti d b - L C z Q 4 @ S C [ i [ S q - > 3 i † ~ > ? i ,, - q @ C C - q C % r i a < - S p  
, i ; b - q f S Y C - ^ @ ^ i 3 C ^ L s > K C ^ C q z S f C - @ f C p - q S Y ^ C z . b q W > | C E j i @ ^ Y S Ci , f- S Y 4 Y C =  
P z e s = w w q f S f i b q L w 4 s w c J C E i | v v c

9J: r<Pqz..SsCq, i Ti>Oi R>- ^@yi Cz - Y> [ -szCq\$^L -z- q\$>Lb><PCss - ^@sPbLS4%eY^^\$^L ..\$P -  
Yc q^C@ \ b@C\$ / -z-q> fbYIDD>eei vE vE>cE| CE

9I: „ i ^ C>ri XS>yi V~q-z <P>di , 44CY- ^@ ^ i K-b> [ -szCq\$^L -z- q\$^L \ Gs ..\$P Y\$ S@ @-z->  
; bpp>fbY -4sw|ccciCEcE| CEci a ^S^Ci , f- S\$ 4Y= Pzes=ww qf.SfibqLw 4sw|ccciCEcE

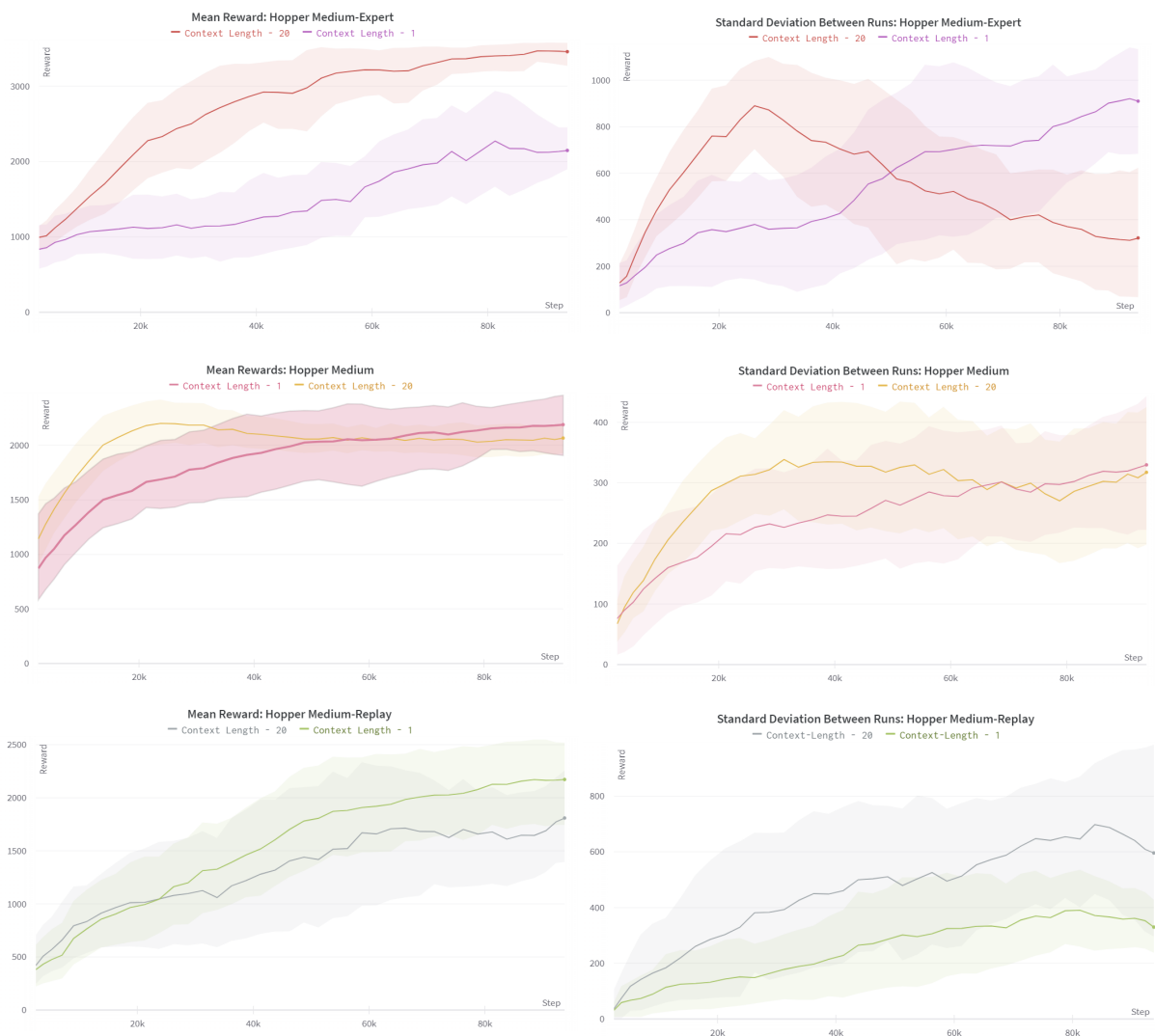




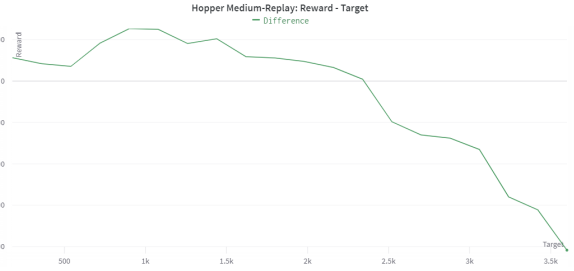
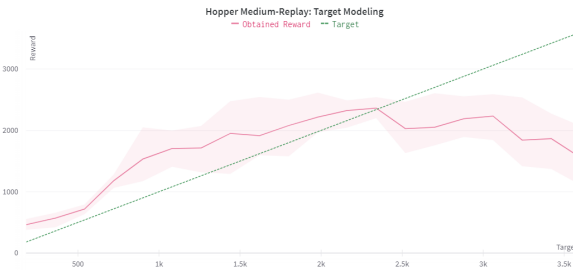
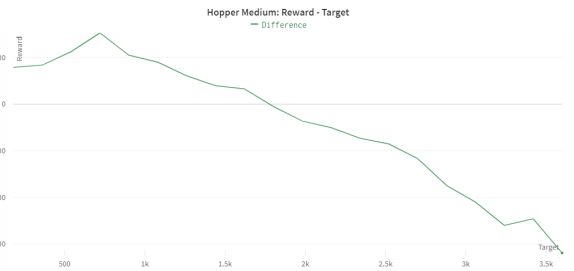
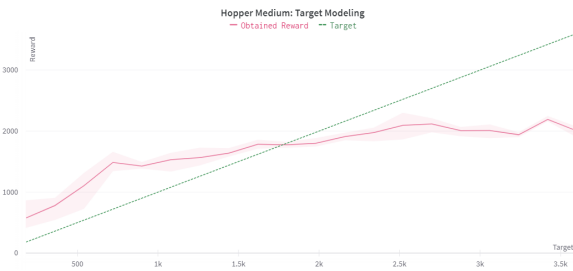
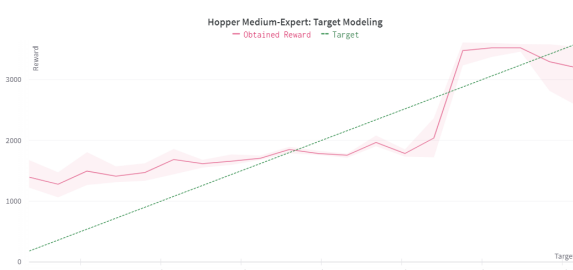


**Decision Transformer validation  
results**

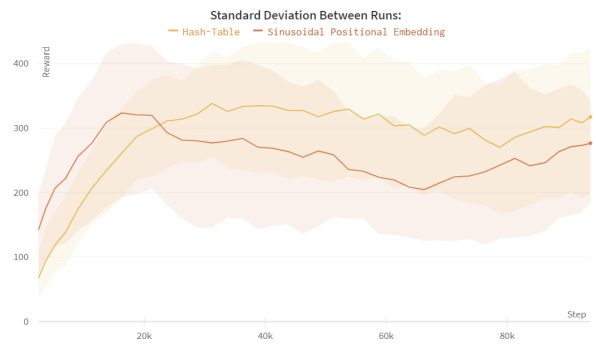
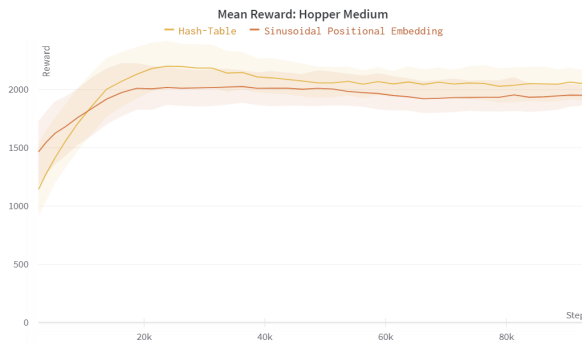
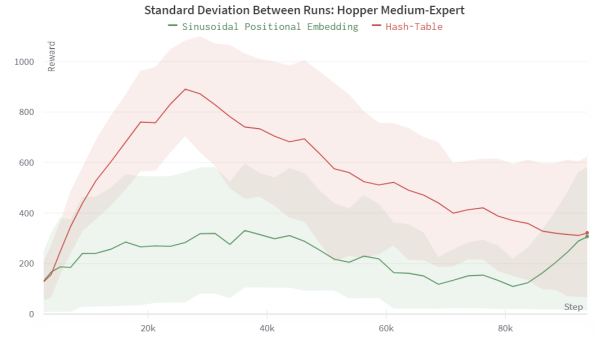
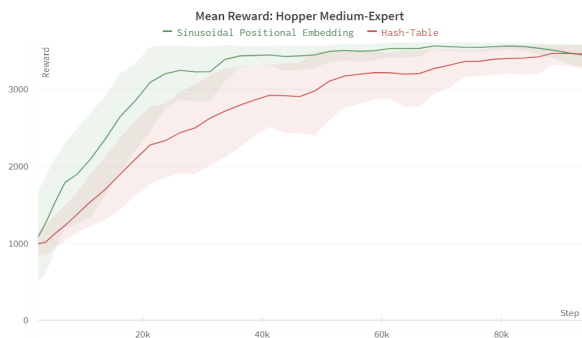
# A.1 Hopper Environment Results



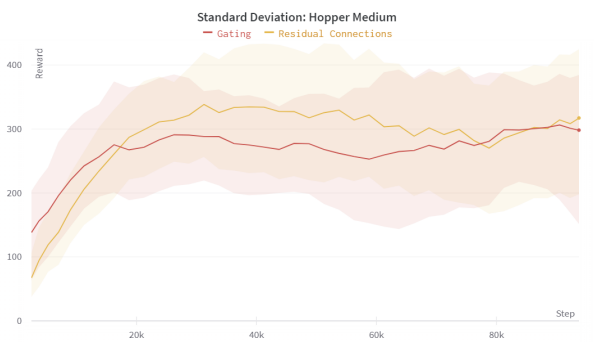
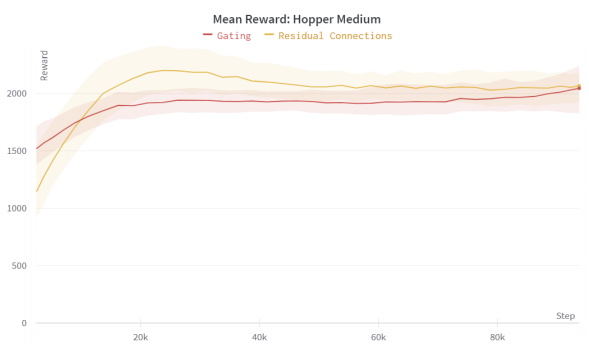
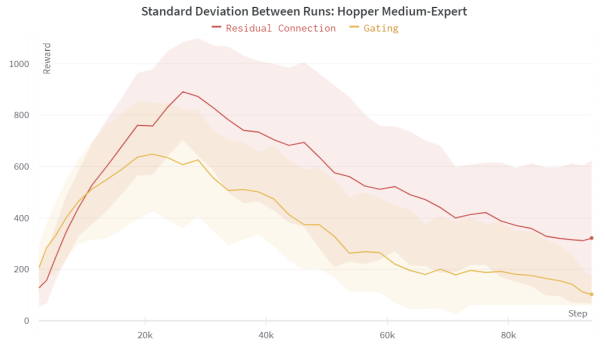
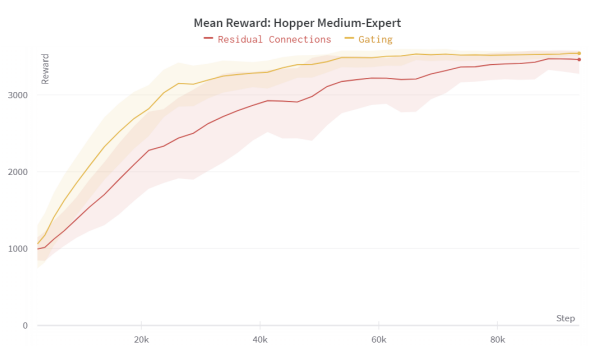
GSL~qC, ic= ; b^zCz XC'LzP= [ G ^ pC.: q@- ^@rz ^@ q@? CjS zB^ - <qss ObeeCq? - z sCz



$G_{i|j} = y_i - q_{\pi}(z_i) + b_{\pi}(z_i) - \langle \phi_{ss}(z_i), \theta_{\pi} \rangle - z_i \cdot c_{\pi}$

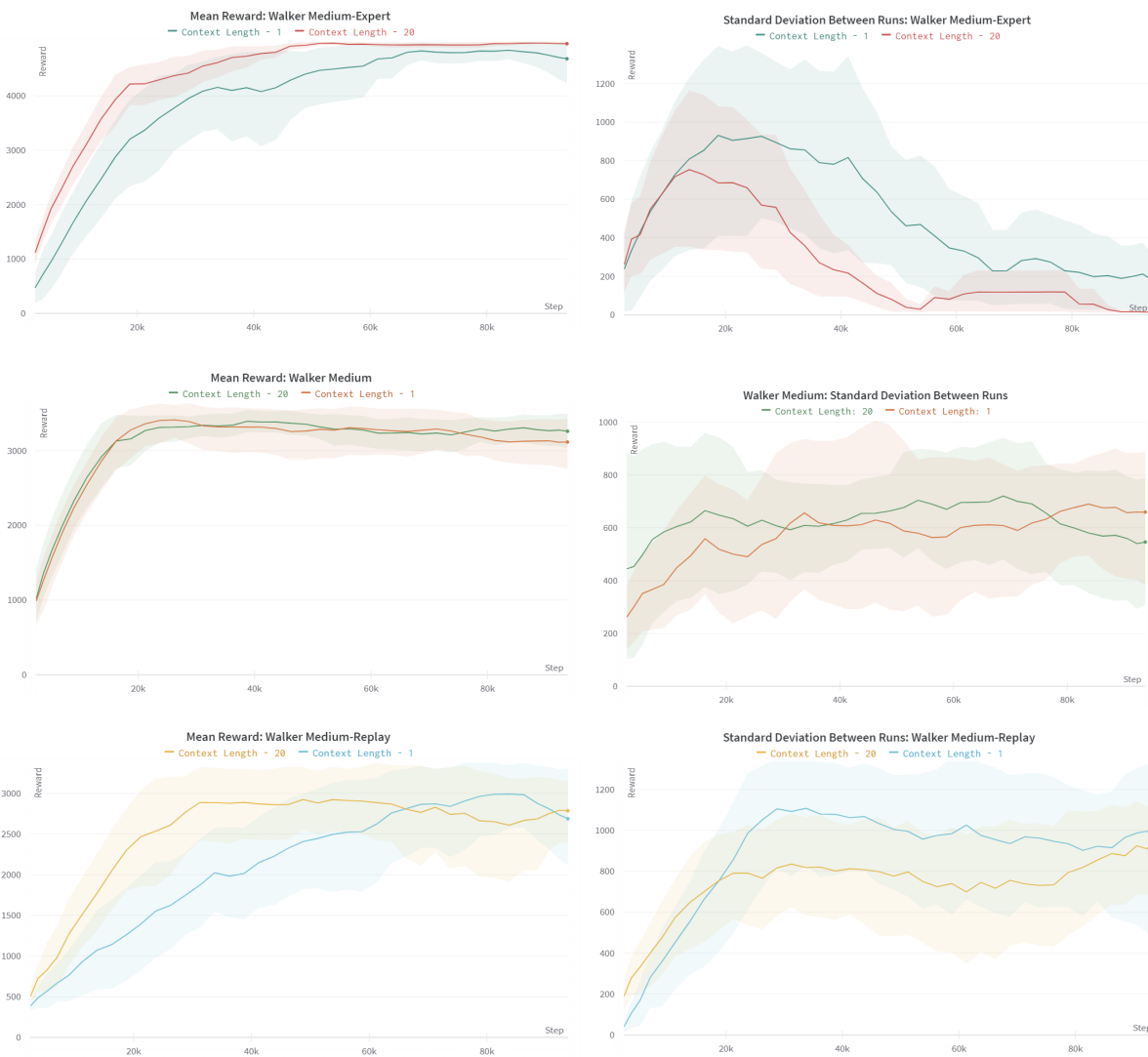


GL-qC, i{= dbSS^ - YB\ 4C@S'L= [ C ^ pC.: q@- ^@rz ^@ q? Cfs zB^ - < qss ObeeCq? - z sCzs

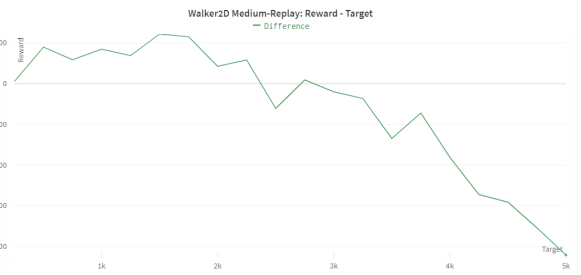
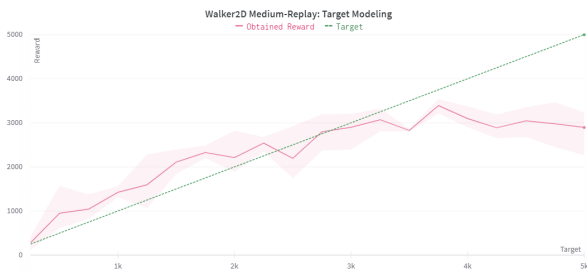
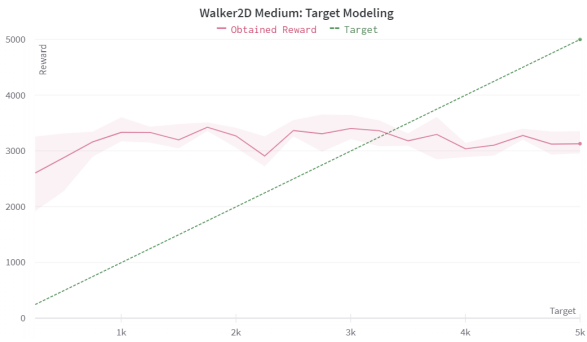
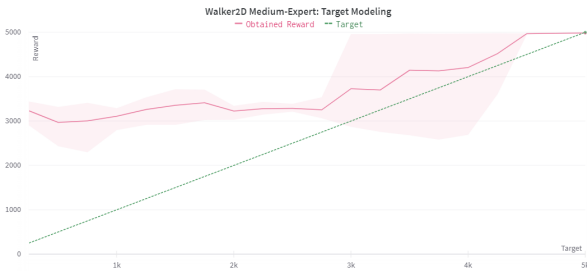


GSL-qC, iJ=K-zSL= [ C ^ pC.: q@- ^@rz ^@ q@? CJS zB^ - < qss ObeeCq? - z sCzs

## A.2 Walker2d Environment Results

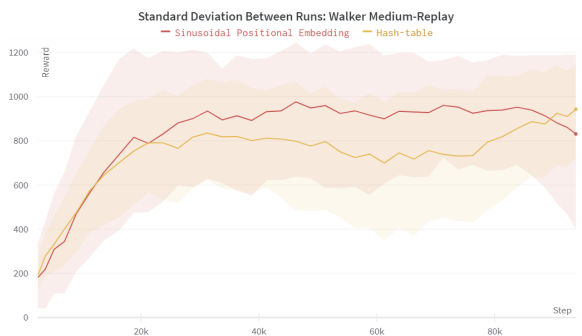
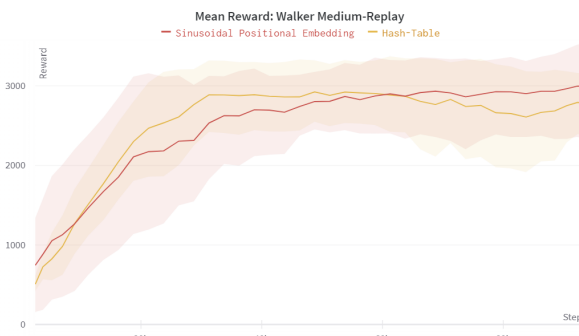
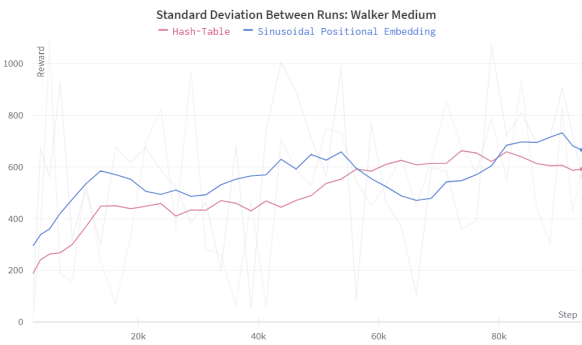
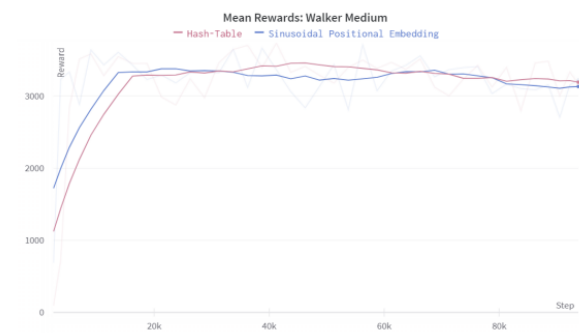
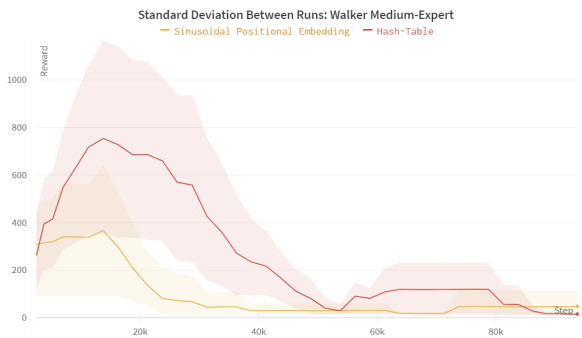
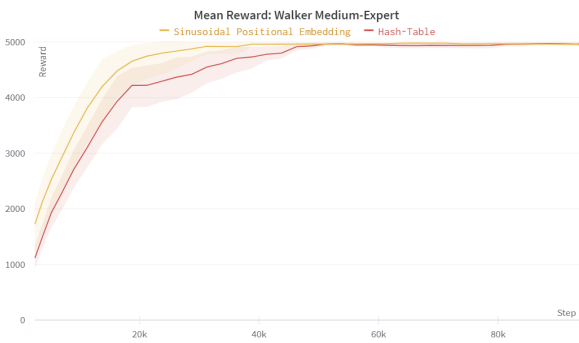


GSL~qC, ii = ; b^zCiz XC^LzP= [ G ^ pC.: q@ - ^@rz ^@ q@? CfS zB^ - <qss ,, - YWq| @? - z- sCzs

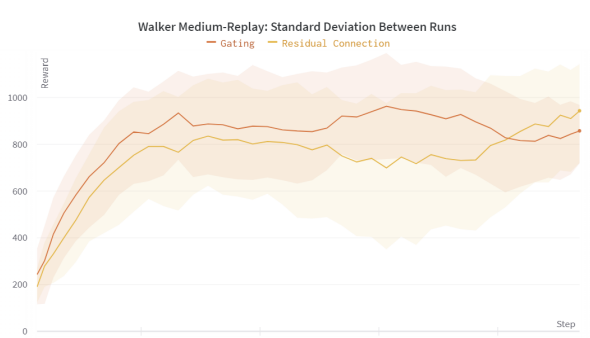
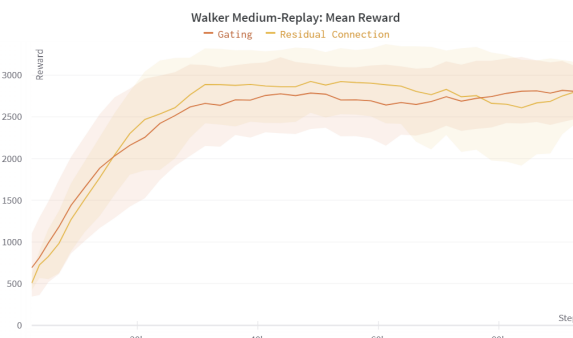
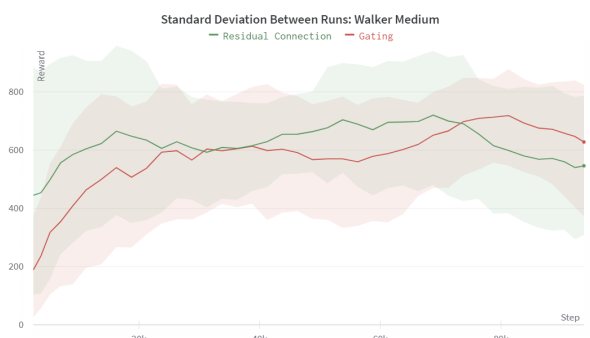
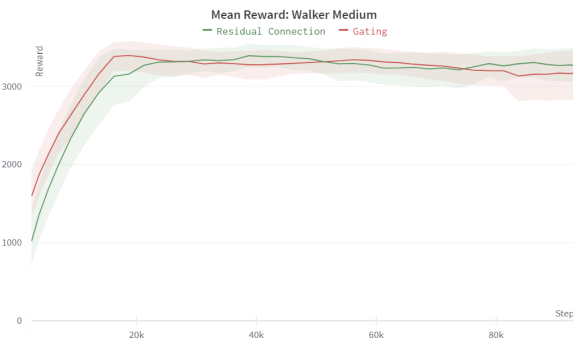
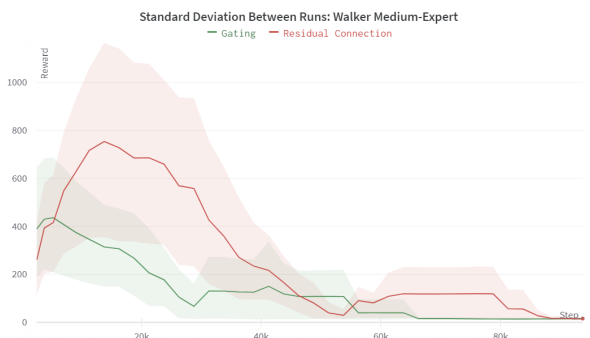
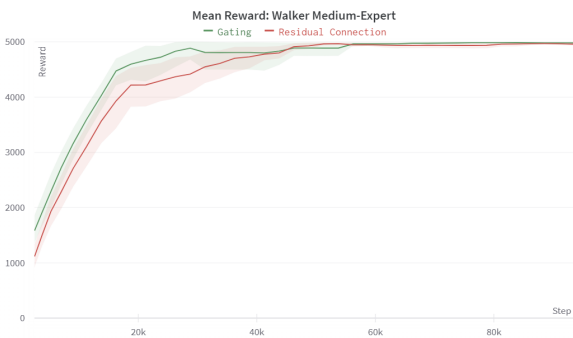


GS~qC, iv= y- qLz [ b@CS'L - <qss „ - Wq? - z sCs





GSL~qC, iu= dbS^B^ - YB\ 4C@8L= [ G ^ pC.: q@ - ^@rz- ^@ q@? CjS zS^ - < qss ,, - Wqj @? - z- sCzs



GSL~qC, iD=K-zSL=[G^pC.:q@-^@rz^@q@?CfSzb^<qpss,, -WUq]@?-z-sCzs