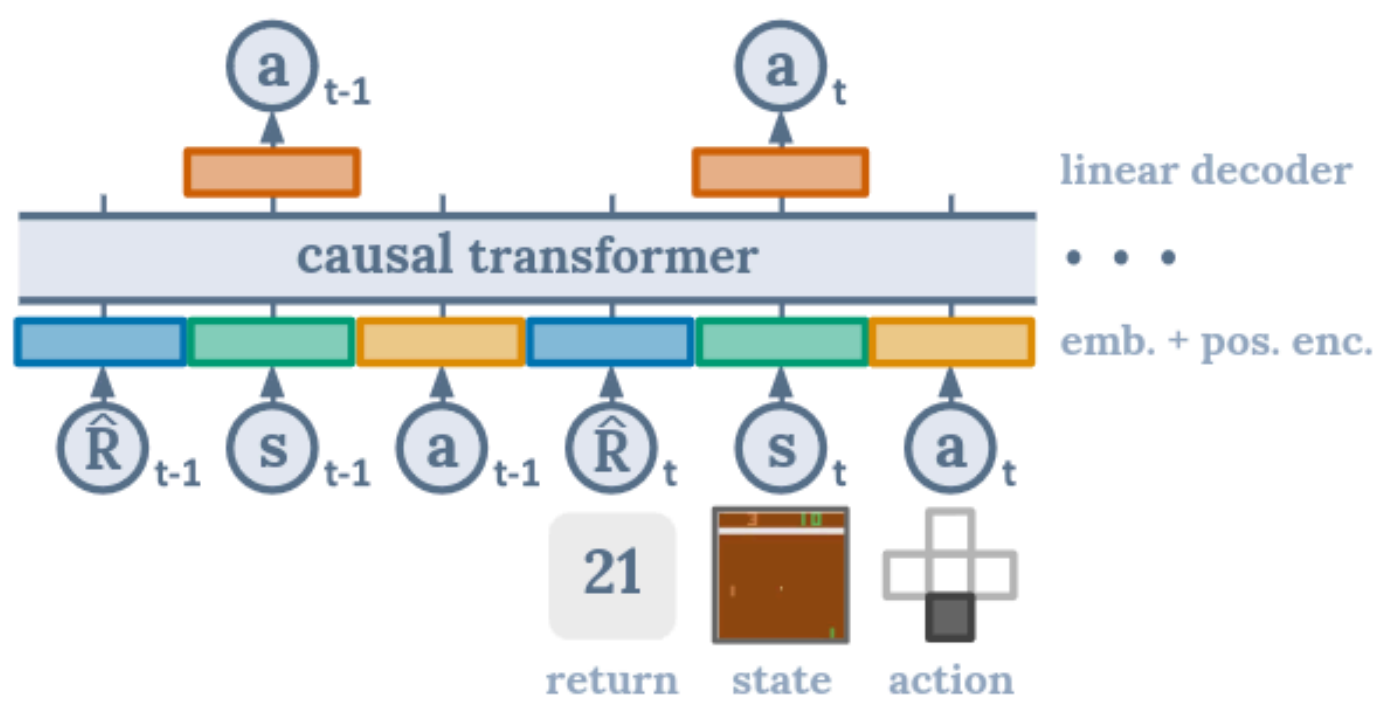


DPT: Pretrained Decision Transformers



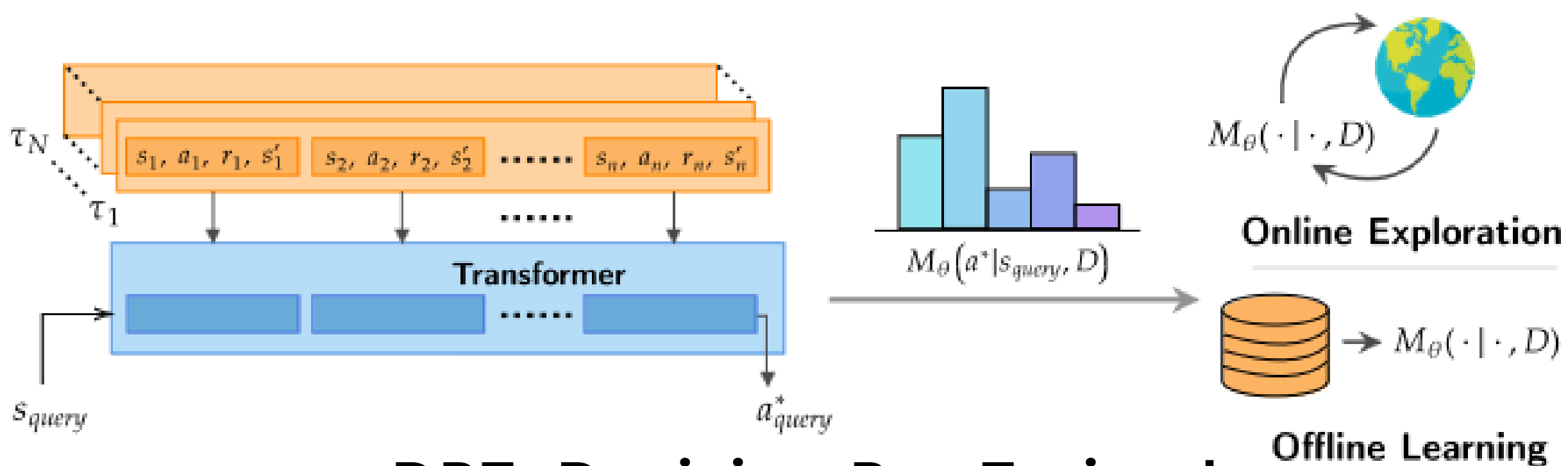
Decision Transformers

- لا بغيثا ندوبو عل كل مفهوم بوحدو كاين:
 - التدريب المُسبق Pretraining: الهندسة المعمارية لـ Transformer ضمن DPT تعالج هذه الثلاثيات الحالة-العمل-المكافأة بكفاءة. هذا التدريب المسبق يعطي النموذج فهما أساسيا للهياكل المكافأة واختيار الإجراءات الأمثل عبر مختلف السيناريوهات.
 - المحوَّلات القرارية (Decision Transformers): على عكس الأساليب التقليدية لتعلم التعزيز اللي كيتركزو بشكل حصري على الحالة الحالية (على سبيل المثال، موقع اللاعبين، مواقع الأعداء)، DPT كيستعمل Decision Transformers. هاد النوع من Transformers مصمم خصيصا باش يأخذ بعين الاعتبار ليس فقط الحالة الحالية بل كذلك تسلسل الحالات والإجراءات السابقة (السياق Context). هاد المعلومات السياقية مهمة لاتخاذ القرارات في المستقبل.

المرة لي دازت هضرنا على ألعاب و كيفاش عووعلي قدرات تخلي المستعمل يخلق لعبة ديالو من صورة او رسمة، دابا نخل وأنت كتدرب برنامج حاسوبي باش يلعب مثلا لعبة فيديو، لكن بدل يد التحكم او clavier، كستعمل وكيل تعلم التعزيز (Reinforcement Learning) لي كيستعمل أحدث التقنيات. الوكلاء التقليدية لتعلم التعزيز كيتعلمو عن طريق التجربة والخطأ، كيبحثو في البيئة ويتلقاو المكافآت على الأفعال الجيدة (التعزيز الإيجابي) والعقوبات على الأفعال السيئة (التعزيز السلبي). هاد النهج يكون بطيء وكيطلب الكثير من البيانات او data، خاصة مع الألعاب المعقدة لي فيها فضاءات حالة عالية الأبعاد.

بوجود التقنيات الجديدة، كيوجد حلا ممكننا. في السنوات الأخيرة، عرفنا تقدما هائلا في مجال التعلم الآلي، خاصة في مجال التعلم بالتعزيز (Reinforcement Learning). ما يميز النهج الجديد، المسمى بـ "Decision-Pretrained Transformer (DPT)"، هو قدرته على تطوير استراتيجيات قرار متعمقة بدون الحاجة إلى كميات ضخمة من البيانات. هاد الشي يعني أن البرنامج قادر على تعلم اللعبة بسرعة أكبر ودقة أكثر، حتى مع الألعاب المعقدة والمتطلبات العالية

هاد النهج الجديد لي سميتو (DPT). DPT كيستفاد من قوة supervised learning و Transformers، هو نوع معيّن من Neural Networks. خلال **التدريب المُشرف** (supervised) نموذج DPT كيتعرض لمجموعة ضخمة من البيانات : الحالة-العمل-المكافأة (State, Action, Reward) من مهام RL متنوعة. هاد لتنوع في المجموعة كيسمح للنموذج يتعلم استراتيجيات و قرارات قابلة للتعميم على بيئات مختلفة.



DPT: Decision Pre-Trained Transformer

- التكيف في السياق In Context Adaptation:
 - منين كيواجه مهمة جديدة في التعلم بالتعزيز، DPT يستفيد من معرفته المسبقة والسياق (التفاعلات السابقة) باش يتنبأ بالإجراء الأمثل في تلك الوضعية المعينة. هذا يسمح لـ DPT بالتكيف بسرعة مع البيئات غير المعروفة وتحقيق كفاءة الحاجة لي كتطلب كمية أقل من data مقارنة بالأساليب التقليدية ديال RL.

نظريا تتعمق الدراسة في الخصائص التقنية لـ DPT عبر استكشاف العلاقة بينه وبين Posterior Sampling، وهي طريقة معروفة بكفاءتها في Reinforcement Learning من زمان. Posterior Sampling ترتبط بقدرة النموذج على استنتاج توزيعات Distributions الاحتمالات الممكنة للنتائج المستقبلية بناءً على المعرفة السابقة والبيانات المتاحة. يعتبر هذا الاتصال مهما بزاف، حيث كيوفر لنموذج DPT واحد أساسا نظري لقدرة على Adaptation مع البيئات المختلفة واستخدام البيانات بكفاءة أكبر.

هاد الارتباط بين DPT وعمليات Bayesian Framework كيساهم في توجيه البحث ل فهم أعمق لكيفية استعمال البيانات بكفاءة فعمليات التعلم، وكيفاش يمكن نحسنو أداء النماذج Models في نفس الوقت. هاد الاتصال كيوجه الابتكارات التقنية للي تطور نماذج أقوى وأكثر كفاءة في ميدان التعلم الآلي و RL.

هاد الطريقة كتحاول تقويم أداء النموذج على أساس قدرته على اتخاذ القرارات الصحيحة في سياقات جديدة ومتغيرة. بمعنى آخر، DPT ماكتكتفيش بمجرد تطبيق ما دربات عليه، بل كتعتمد على هاد **transformers** لتوجيه القرارات الأمثل في الظروف الجديدة. هاد الشي كيجعل النموذج أكثر قابلية للتكيف والتطبيق في مجموعة متنوعة من الألعاب والمهام، مما كيزيد من أدائه وفاعليته.

المفهوم الأساسي هنا هو **السياق** (Context) او In-Context Learning ICL. السياق يشير لي اللي حصل قبل فاللعبة – الحركات لي درت، والعراقيل لي واجهت، والمكافآت لي حصلت عليها. بواسطة تذكر هادشي، DPT قادر يحسن من قراراته فالوضعية الحالية.

على سبيل المثال، إذا كنت مطارّد من طرف عدو في اطار لعبة، DPT قادر يستغل معرفتو باستراتيجيات الهروب (تعلمها من ألعاب أخرى) باش يتجنب الخطر.

هاد التجميع ديال التدريب المسبق، Decision Transformers، والوعي ب Context كيمكن لـ DPT يتعلم الألعاب الجديدة بشكل أسرع بزاف من الوكلاء التقليديين لـ RL. هو ماكيحتاجش لتمرين كثير وكيقدر يتكيف استراتيجيتو بناء على اللعبة.

واحد ملاحظة من الضروري الاعتراف بأن DPT هو مجال بحث نسبيًا جديد حيث عاد خرج البحث العام لي داز، ومن الضروري إجراء مزيد من البحوث لفهم فعاليته بشكل كامل عبر مشاكل ديال Reinforcement Learning.

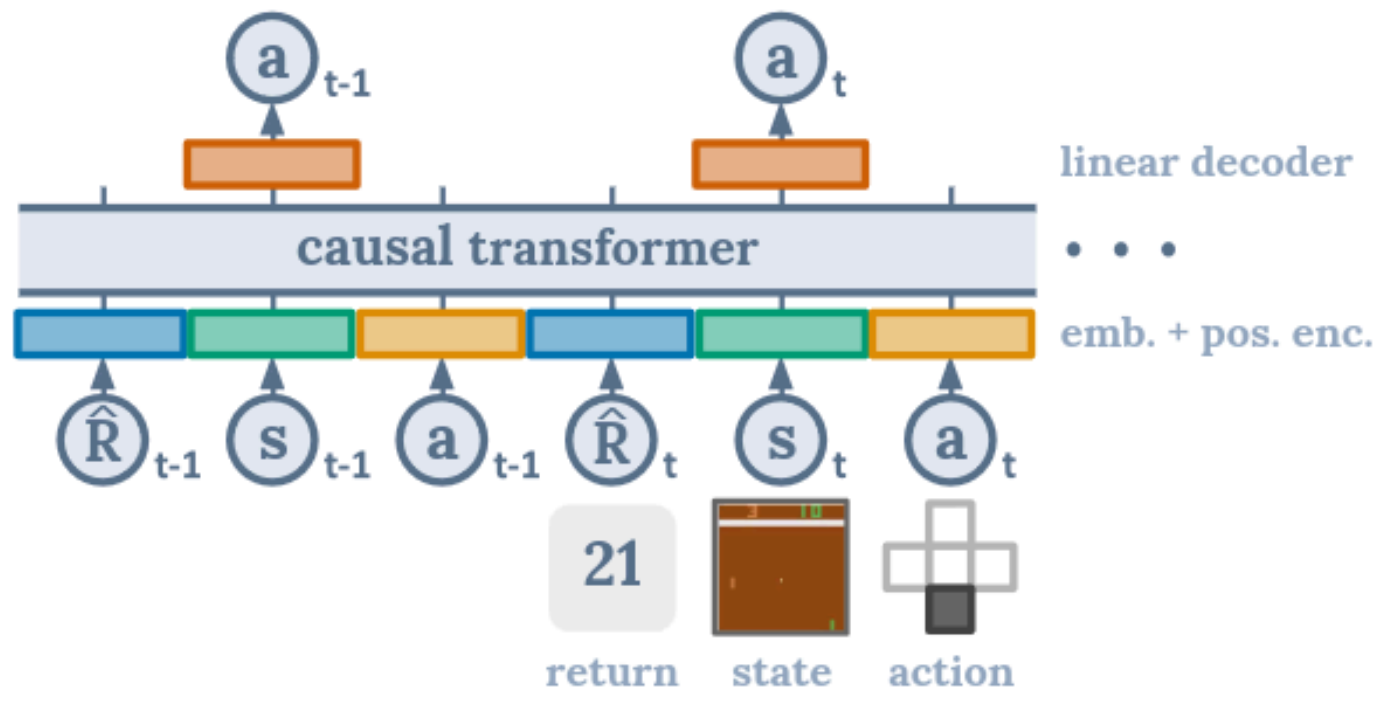
References:

- Decision Transformer: Reinforcement Learning via Sequence Modeling: <https://arxiv.org/pdf/2106.01345>
- Supervised Pretraining Can Learn In-Context Reinforcement Learning: <https://arxiv.org/pdf/2306.14892>

official partner

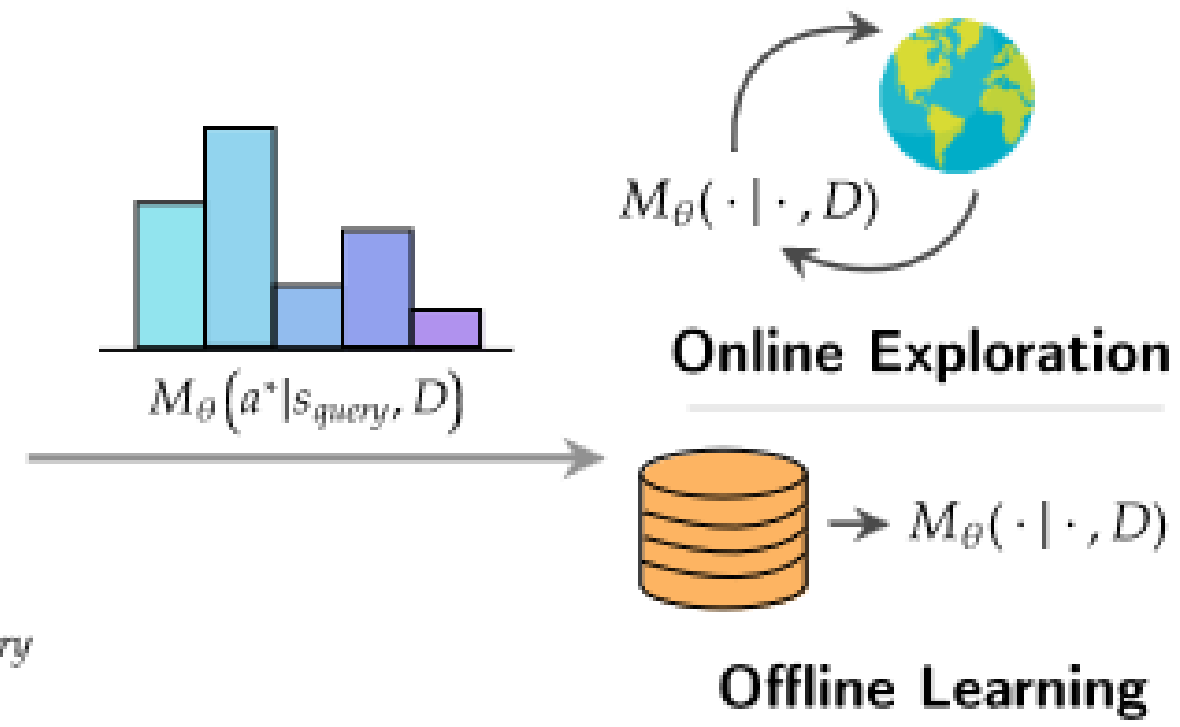
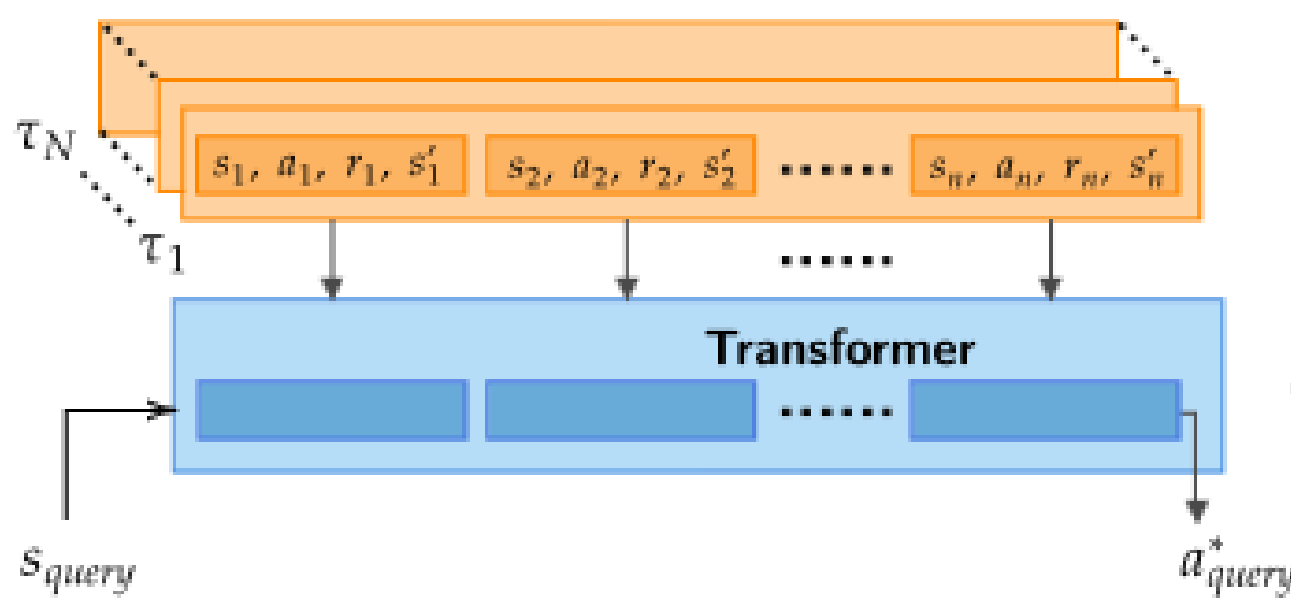


DPT: Pretrained Decision Transformers



Decision Transformers

- إذا أردنا التفصيل في كل مفهوم على حدة:
 - **التدريب المُسبق (Pretraining):** هندسة المحولات (Transformers) ضمن نموذج DPT تمكن من معالجة هذه الثلاثيات الحالة-العمل-المكافأة (State-Action-Reward) بكفاءة. يعطي هذا التدريب المُسبق النموذج فهماً أساسياً لهياكل المكافأة واختيار الإجراءات الأمثل عبر مختلف السيناريوهات.
 - **المحولات القرارية (Decision Transformers):** على عكس الأساليب التقليدية لتعلم التعزيز التي تركز بشكل حصري على الحالة الحالية (على سبيل المثال، موقع اللاعب، مواقع الأعداء)، يستخدم DPT المحولات القرارية. هذا النوع من المحولات مصمم خصيصاً ليأخذ بعين الاعتبار ليس فقط الحالة الحالية بل أيضاً تسلسل الحالات والإجراءات السابقة (السياق). هذه المعلومات السياقية مهمة لاتخاذ القرارات في المستقبل



DPT: Decision Pre-Trained Transformers

تكيف السياق In Context Adaptation: عندما يواجه نموذج DPT مهمة جديدة في التعلم بالتعزيز، يستفيد من معرفته المسبقة والسياق (التفاعلات السابقة) ليتنبأ بالإجراء الأمثل في تلك الحالة المحددة. يتيح ذلك لـ DPT التكيف بسرعة مع البيئات غير المألوفة وتحقيق كفاءة باستخدام كمية أقل من البيانات مقارنة بالأساليب التقليدية في تعلم التعزيز.

من الناحية النظرية، يشمل التعمق في الدراسة للخصائص التقنية لـ **DPT** استكشاف العلاقة بينه وبين **Posterior Sampling**، وهي طريقة معروفة بكفاءتها في التعلم بالتعزيز منذ فترة طويلة. ترتبط Posterior Sampling بقدرة النموذج على استنتاج توزيعات احتمالية لنتائج المستقبل بناءً على المعرفة السابقة والبيانات المتاحة. يُعتبر هذا الارتباط مهماً لأنه يوفر لـ DPT أساساً نظرياً لتكيفه مع بيانات متنوعة واستخدام البيانات بكفاءة أكبر.

يسهم هذا الاتصال بين DPT و Bayesian posterior sampling في توجيه البحث نحو فهم أعمق لاستخدام البيانات بكفاءة في عمليات التعلم، وكيفية تحسين أداء النماذج في نفس الوقت. يوجه هذا الارتباط الابتكارات التقنية التي تُطوّر نماذجاً أقوى وأكثر كفاءة في مجالات التعلم الآلي والتعلم بالتعزيز (RL).

هذه الطريقة تحاول تقييم أداء النموذج استناداً إلى قدرته على اتخاذ القرارات الصحيحة في سياقات جديدة ومتغيرة. بمعنى آخر، فإن DPT لا يكتفي بمجرد تطبيق ما تم تدريبه عليه، بل يعتمد على هذه المحولات القرارية لتوجيه القرارات الأمثل في الظروف الجديدة.

هذا الأسلوب يجعل النموذج أكثر قابلية للتكيف المرنة في مجموعة متنوعة من الألعاب والمهام، مما يزيد من أدائه وفاعليته. فيفضل استخدام المحولات القرارية، يمكن لـ DPT تحديث وتعديل استراتيجياته واتخاذ القرارات الأمثل بناءً على السياقات الجديدة التي يواجهها، مما يعزز قدرته على التكيف وتحسين أدائه في بيئات مختلفة ومتغيرة.

المفهوم الأساسي هنا هو السياق (Context) أو التعلم في السياق (In-Context Learning - ICL). يشير السياق إلى كل ما حدث مسبقاً في اللعبة، مثل الحركات التي قمت بها، والعقبات التي واجهتها، والمكافآت التي حصلت عليها. من خلال تذكر هذه الأحداث، يمكن لـ DPT تحسين قراراته في الوضعية الحالية.

على سبيل المثال، إذا كنت تُطارَد من قِبَل عدو، يمكن لـ DPT استخدام معرفته باستراتيجيات الهروب التي تعلمها من ألعاب أخرى لتجنب الخطر.

هذا التجميع للتدريب المسبق، المحولات القرارية، والوعي بالسياق يمكن أن يساعد DPT على تعلم الألعاب الجديدة بشكل أسرع بكثير من الوكلاء التقليديين لتعلم التعزيز. فعندما يتم تدريب DPT بشكل جيد مسبقاً، يكون لديه فهم أساسي للهيكل والأنماط المشتركة في ألعاب التعلم بالتعزيز.

المصادر :

- Decision Transformer: Reinforcement Learning via Sequence Modeling: <https://arxiv.org/pdf/2106.01345>
- Supervised Pretraining Can Learn In-Context Reinforcement Learning: <https://arxiv.org/pdf/2306.14892>

official partner



