

# VAI TRÒ THIỆT YẾU CỦA “RÀO CHẮN” TRONG HỆ THỐNG AI



➡ THẢO LÂM

Thuật ngữ “rào chắn” thường đề cập đến các cơ chế, chính sách và thực tiễn được đưa ra để đảm bảo rằng hệ thống trí tuệ nhân tạo (AI) hoạt động an toàn và trong các ranh giới nhất định được xác định trước, đặc biệt là với các mô hình AI tổng hợp. Những rào chắn này rất quan trọng để đảm bảo việc sử dụng AI có trách nhiệm và giảm thiểu rủi ro cũng như những hậu quả không lường trước được.

**T**rí tuệ nhân tạo (AI) đã thâm nhập vào cuộc sống hàng ngày của chúng ta, trở thành một phần không thể thiếu trong nhiều lĩnh vực khác nhau - từ chăm sóc sức khỏe và giáo dục đến giải trí và tài chính. Công nghệ đang phát triển với tốc độ chóng mặt, giúp cuộc sống của chúng ta dễ dàng hơn, hiệu quả hơn và thú vị hơn theo nhiều cách. Tuy nhiên, giống như bất kỳ công cụ mạnh mẽ nào khác, AI cũng tiềm ẩn những rủi ro có hậu, đặc biệt khi được sử dụng một cách vô trách nhiệm hoặc không có sự giám sát đầy đủ.

Điều này dẫn đến một thành phần thiết yếu của hệ thống AI - rào chắn. Các rào chắn trong hệ thống AI đóng vai trò là biện pháp bảo vệ để đảm bảo việc sử dụng công nghệ AI có trách nhiệm. Chúng bao gồm các chiến lược, cơ chế và chính sách được thiết kế để ngăn chặn việc lạm dụng, bảo vệ quyền riêng tư của người dùng và thúc đẩy tính minh bạch và công bằng.

## **Rào chắn trong hệ thống AI là gì?**

Các công nghệ AI, do tính chất tự chủ và thường tự học, đặt ra những thách thức đặc biệt. Những thách thức này đòi hỏi phải có một bộ nguyên tắc hướng dẫn và kiểm soát cụ thể - các rào chắn. Chúng rất cần thiết trong việc thiết kế và triển khai hệ thống AI, xác định ranh giới của hành vi AI có thể chấp nhận được.

Rào chắn trong hệ thống AI bao gồm nhiều khía cạnh. Chủ yếu, chúng phục vụ để bảo vệ chống lại việc lạm dụng, thiên vị và các hành vi phi đạo đức. Điều này bao gồm việc đảm bảo các công nghệ AI hoạt động theo các thông số đạo đức do xã hội đặt ra và tôn trọng quyền riêng tư cũng như quyền của cá nhân.

Các rào chắn trong hệ thống AI có nhiều dạng khác nhau, tùy thuộc vào đặc điểm cụ thể của hệ thống AI và mục đích sử dụng của nó. Ví dụ: Chúng có thể bao gồm các cơ chế đảm bảo quyền riêng tư và bảo mật dữ liệu, các quy trình ngăn

chặn kết quả phân biệt đối xử và các chính sách bắt buộc phải kiểm tra thường xuyên các hệ thống AI để tuân thủ các tiêu chuẩn đạo đức và pháp lý.

Một phần quan trọng khác của rào chắn là tính minh bạch - đảm bảo rằng các quyết định do hệ thống AI đưa ra có thể được hiểu và giải thích được. Tính minh bạch cho phép trách nhiệm giải trình, đảm bảo rằng các lỗi hoặc hành vi sử dụng sai có thể được xác định và khắc phục.

Hơn nữa, các rào chắn có thể bao gồm các chính sách bắt buộc phải có sự giám sát của con người trong các quá trình ra quyết định quan trọng. Điều này đặc biệt quan trọng trong các tình huống rủi ro cao, trong đó các lỗi AI có thể dẫn đến tác hại đáng kể, chẳng hạn như trong lĩnh vực chăm sóc sức khỏe hoặc xe tự hành.

Cuối cùng, mục đích của các rào chắn trong hệ thống AI là để đảm bảo rằng công nghệ AI giúp nâng cao khả năng của con người và làm phong phú thêm cuộc sống của chúng ta mà không ảnh hưởng đến các quyền và sự an toàn của chúng ta. Chúng đóng vai trò là cầu nối giữa tiềm năng to lớn của AI và khả năng hiện thực hóa nó một cách an toàn và có trách nhiệm.

## **Tầm quan trọng của rào chắn trong hệ thống AI**

Trong viễn cảnh năng động của công nghệ AI, tầm quan trọng của các rào chắn không thể bị phóng đại. Khi các hệ thống AI phát triển phức tạp và tự chủ hơn, chúng được giao những nhiệm vụ có tác động và trách nhiệm lớn hơn. Do đó, việc triển khai hiệu quả các rào chắn không chỉ mang lại lợi ích mà còn cần thiết để AI phát huy hết tiềm năng một cách có trách nhiệm.

Lý do đầu tiên giải thích tầm quan trọng của rào chắn trong hệ thống AI nằm ở khả năng bảo vệ khỏi việc lạm dụng công nghệ AI. Khi các hệ thống AI có được nhiều khả năng hơn, nguy cơ các hệ thống này bị sử dụng cho mục đích xấu sẽ

tăng lên. Rào chắn có thể giúp thực thi các chính sách sử dụng và phát hiện hành vi sử dụng sai mục đích, giúp đảm bảo rằng công nghệ AI được sử dụng một cách có trách nhiệm.

Một khía cạnh quan trọng khác về tầm quan trọng của rào chắn là đảm bảo sự công bằng và chống lại sự thiên vị. Hệ thống AI học hỏi từ dữ liệu mà chúng được cung cấp và nếu dữ liệu này phản ánh những thành kiến xã hội, thì hệ thống AI có thể duy trì và thậm chí khuếch đại những thành kiến này. Bằng cách triển khai các rào chắn để tích cực tìm kiếm và giảm thiểu những thành kiến trong việc ra quyết định bằng AI, chúng ta có thể đạt được những bước tiến hướng tới các hệ thống AI công bằng hơn.

Các rào chắn cũng rất cần thiết trong việc duy trì niềm tin của công chúng vào công nghệ AI. Tính minh bạch, được hỗ trợ bởi các rào chắn, giúp đảm bảo rằng các quyết định do hệ thống AI đưa ra có thể được hiểu và thẩm vấn. Sự cởi mở này không chỉ thúc đẩy trách nhiệm giải trình mà còn góp phần tạo dựng niềm tin của công chúng vào công nghệ AI.

Hơn nữa, rào chắn rất quan trọng để tuân thủ các tiêu chuẩn pháp lý và quy định. Khi các chính phủ và cơ quan quản lý trên toàn thế giới nhận ra những tác động tiềm ẩn của AI, họ đang thiết lập các quy định để quản lý việc sử dụng AI. Việc triển khai hiệu quả các rào chắn có thể giúp các hệ thống AI luôn nằm trong các ranh giới pháp lý này, giảm thiểu rủi ro và đảm bảo hoạt động trơn tru.

Các rào chắn cũng tạo điều kiện thuận lợi cho sự giám sát của con người trong các hệ thống AI, củng cố khái niệm về AI như một công cụ hỗ trợ chứ không phải thay thế việc ra quyết định của con người. Bằng cách thực hiện con người trong vòng lặp HITL, đặc biệt là trong các quyết định mang tính rủi ro cao, các rào chắn có thể giúp đảm bảo rằng các hệ thống AI vẫn nằm trong tầm kiểm soát của chúng ta và các quyết định

của chúng phù hợp với các giá trị và chuẩn mực chung của chúng ta.

Về bản chất, việc triển khai các rào chắn trong hệ thống AI là hết sức quan trọng để khai thác sức mạnh biến đổi của AI một cách có trách nhiệm. Chúng đóng vai trò là bức tường thành chống lại những rủi ro và cạm bẫy tiềm ẩn liên quan đến việc triển khai các công nghệ AI, khiến chúng trở thành một phần không thể thiếu trong tương lai của AI.

### Sự trỗi dậy của AI tạo sinh

Sự ra đời của các hệ thống AI tạo sinh như ChatGPT của OpenAI và Bard của Google đã nhấn mạnh hơn nữa sự cần thiết của các biện pháp bảo vệ mạnh mẽ trong các hệ thống AI. Những mô hình ngôn ngữ phức tạp này có khả năng tạo ra văn bản giống con người, tạo ra phản hồi, câu chuyện hoặc bài viết kỹ thuật chỉ trong vài giây. Khả năng này tuy ấn tượng và vô cùng hữu ích nhưng cũng tiềm ẩn những rủi ro.

Hệ thống AI tạo sinh có thể tạo ra nội dung có thể không phù hợp, có hại hoặc lừa đảo nếu không được giám sát đầy đủ. Họ có thể tuyên truyền những thành kiến được nhúng trong dữ liệu đào tạo của họ, có khả năng dẫn đến kết quả đầu ra phản ánh quan điểm phân biệt đối xử hoặc thành kiến. Ví dụ: nếu không có rào chắn bảo vệ thích hợp, những mô hình này có thể được sử dụng để tạo ra thông tin sai lệch hoặc tuyên truyền có hại.

Hơn nữa, các khả năng tiên tiến của AI tạo ra cũng giúp tạo ra thông tin thực tế nhưng hoàn toàn hư cấu. Nếu không có các biện pháp bảo vệ hiệu quả, điều này có thể bị sử dụng với mục đích xấu nhằm tạo ra những tường thuật sai sự thật hoặc truyền bá thông tin sai lệch. Quy mô và tốc độ hoạt động của các hệ thống AI này sẽ làm tăng thêm tác hại tiềm tàng của việc sử dụng sai mục đích đó.

Do đó, với sự gia tăng của các hệ thống AI tạo sinh mạnh mẽ, nhu cầu về rào chắn lại càng quan trọng hơn nữa. Chúng giúp đảm bảo những công nghệ này được sử dụng một cách có trách nhiệm và có đạo đức, thúc đẩy tính minh bạch, trách nhiệm giải trình và tôn trọng các chuẩn mực và giá trị xã hội. Về bản chất, các rào chắn bảo vệ chống lại việc lạm dụng AI, đảm bảo khả năng tạo ra tác động tích cực của nó đồng thời giảm thiểu nguy cơ gây hại.

### Những thách thức và giải pháp

Triển khai các rào chắn trong hệ thống AI là một quá trình phức tạp, đặc biệt là vì những thách thức kỹ thuật liên quan. Tuy nhiên, những điều này không phải là không thể vượt qua và có một số chiến lược mà các tổ chức, doanh nghiệp có thể sử dụng để đảm bảo hệ thống AI của họ hoạt động trong giới hạn được xác định trước.

#### *Những thách thức và giải pháp kỹ thuật*

Nhiệm vụ áp đặt các rào chắn trên hệ thống AI thường liên quan đến việc điều hướng một mê cung phức tạp về mặt kỹ thuật. Tuy nhiên, các tổ chức, doanh nghiệp có thể áp dụng cách tiếp cận chủ động bằng cách sử dụng các kỹ thuật máy học mạnh mẽ, như đào tạo đối nghịch và quyền riêng tư khác biệt.

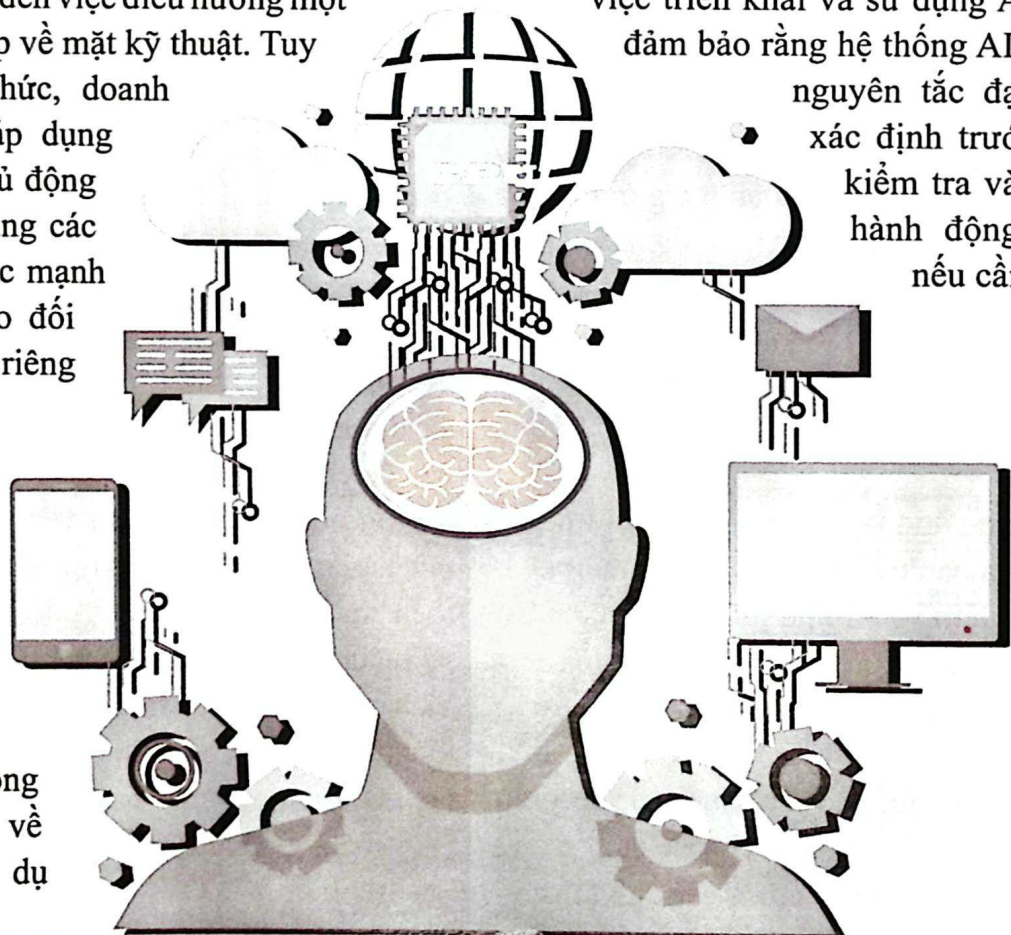
*Đào tạo đối nghịch* là một quá trình bao gồm việc đào tạo mô hình AI không chỉ về các thông tin đầu vào mong muốn mà còn về một loạt các ví dụ

đối nghịch được tạo ra. Những ví dụ đối nghịch này là các phiên bản đã được điều chỉnh của dữ liệu gốc, nhằm mục đích đánh lừa mô hình mắc lỗi. Bằng cách học hỏi từ những đầu vào bị thao túng này, hệ thống AI trở nên tốt hơn trong việc chống lại các nỗ lực khai thác các lỗ hổng của nó.

*Quyền riêng tư khác biệt* là phương pháp thêm nhiều vào dữ liệu huấn luyện để che khuất các điểm dữ liệu riêng lẻ, do đó bảo vệ quyền riêng tư của các cá nhân trong tập dữ liệu. Bằng cách đảm bảo quyền riêng tư của dữ liệu đào tạo, các tổ chức, doanh nghiệp có thể ngăn hệ thống AI vô tình học và truyền bá thông tin nhạy cảm.

#### *Những thách thức và giải pháp vận hành*

Ngoài những phức tạp về mặt kỹ thuật, khía cạnh vận hành của việc thiết lập các rào chắn bảo vệ AI cũng có thể gặp nhiều thách thức. Vai trò và trách nhiệm rõ ràng cần được xác định trong một tổ chức để giám sát và quản lý hiệu quả các hệ thống AI. Một hội đồng hoặc ủy ban đạo đức AI có thể được thành lập để giám sát việc triển khai và sử dụng AI. Họ có thể đảm bảo rằng hệ thống AI tuân thủ các nguyên tắc đạo đức được xác định trước, tiến hành kiểm tra và đề xuất các hành động khắc phục nếu cần thiết.



Hơn nữa, các tổ chức, doanh nghiệp cũng nên xem xét triển khai các công cụ để ghi nhật ký và kiểm tra kết quả đầu ra của hệ thống AI cũng như quá trình ra quyết định. Những công cụ như vậy có thể giúp truy tìm nguyên nhân gốc rễ của bất kỳ quyết định gây tranh cãi nào do AI đưa ra, từ đó cho phép sửa chữa và điều chỉnh một cách hiệu quả.

**Những thách thức và giải pháp về pháp lý và quy định**

Sự phát triển nhanh chóng của công nghệ AI thường vượt xa các khuôn khổ pháp lý và quy định hiện có. Do đó, các tổ chức, doanh nghiệp có thể phải đối mặt với sự không chắc chắn về các vấn đề tuân thủ khi triển khai hệ thống AI. Tương tác với các cơ quan pháp lý và quản lý, cập nhật thông tin về các luật AI mới và chủ động áp dụng các phương pháp hay nhất có thể giảm thiểu những lo ngại này. Các tổ chức, doanh nghiệp cũng nên ủng hộ quy định công bằng và hợp lý trong lĩnh vực AI để đảm bảo sự cân bằng giữa đổi mới và an toàn.

Việc triển khai các rào chắn AI không phải là nỗ lực một lần mà đòi hỏi phải theo dõi, đánh giá và điều chỉnh liên tục. Khi công nghệ AI tiếp tục phát triển thì nhu cầu về các chiến lược đổi mới để bảo vệ khỏi việc lạm dụng cũng tăng theo. Bằng cách nhận biết và giải quyết những thách thức liên quan đến việc triển khai các rào chắn bảo vệ AI, các tổ chức, doanh nghiệp có thể đảm bảo tốt hơn việc sử dụng AI có trách nhiệm.

**Tại sao rào chắn AI nên là trọng tâm chính?**

Khi chúng ta tiếp tục vượt qua các giới hạn về những gì AI có thể làm, việc đảm bảo các hệ thống này hoạt động trong giới hạn đạo đức và trách nhiệm ngày càng trở nên quan trọng. Rào chắn đóng một vai trò quan trọng trong việc duy trì sự an toàn, công bằng và minh bạch của hệ thống AI. Chúng hoạt động như những trạm kiểm soát cần thiết nhằm ngăn chặn khả năng

lạm dụng công nghệ AI, đảm bảo rằng chúng ta có thể thu được lợi ích từ những tiến bộ này mà không ảnh hưởng đến các nguyên tắc đạo đức hoặc gây ra những tổn hại ngoài ý muốn.

Việc triển khai các rào chắn AI đặt ra một loạt thách thức về kỹ thuật, vận hành và quy định. Tuy nhiên, thông qua quá trình đào tạo đối nghịch nghiêm ngặt, các kỹ thuật bảo mật khác biệt và việc thành lập các hội đồng đạo đức AI, những thách thức này có thể được giải quyết một cách hiệu quả. Hơn nữa, hệ thống kiểm tra và ghi nhật ký mạnh mẽ có thể giữ cho quá trình ra quyết định của AI trở nên minh bạch và có thể theo dõi được.

Theo các nhà nghiên cứu, thách thức lớn nhất đối với sự an toàn của AI là tìm hiểu xem các rào chắn có thực sự hoạt động hay không. Hiện tại, rất khó để xây dựng cách đánh giá cho các rào chắn bảo vệ AI vì các mô hình kết thúc mở, có thể được hỏi vô số câu hỏi và trả lời theo vô số cách khác nhau. *“Nó giống như việc cố gắng tìm hiểu tính cách của một người bằng cách nói chuyện với họ. Đó chỉ là một nhiệm vụ khó khăn và phức tạp”* - Amodei của Anthropic (Công ty hiện đang nghiên cứu cách sử dụng chính AI để tạo ra những đánh giá).

Trong tương lai, nhu cầu về rào chắn AI sẽ chỉ tăng lên khi chúng ta ngày càng phụ thuộc vào hệ thống AI. Đảm bảo việc sử dụng có trách nhiệm là trách nhiệm chung - một trách nhiệm đòi hỏi nỗ lực phối hợp của các nhà phát triển, người dùng và cơ quan quản lý AI. Bằng cách đầu tư vào việc phát triển và triển khai các rào chắn AI, chúng ta có thể thúc đẩy một bối cảnh công nghệ không chỉ đổi mới mà còn phù hợp và an toàn. ■ TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN

Tài liệu tham khảo:  
 1. <https://www.ft.com/content/>  
 2. <https://www.unite.ai/>  
 3. <https://www.cigionline.org/>

THƯ VIỆN TP. CẦN THƠ