

AI Regulation Blueprint

A high level blueprint for domestic regulation of civilian advanced AI models.

THREAT MODELS ADDRESSED

- Emergence of unexpected, dangerous behaviour
- Malicious use by known or low-resource bad actors
- Societal disruption from surprise release of a powerful model
- Concentration of power by model developers

THREAT MODELS NOT CURRENTLY ADDRESSED

(NOT EXHAUSTIVE)

- National security use of models
- Malicious use by unknown, high-resource bad actors
- Scenarios where deceptive alignment emerges with no prior warning between scaling steps
- Geopolitical conflict due to fear of “falling behind”

CONTENTS

I Before Training

[SKIP TO SECTION I](#) →

II Training for Broad Competence

[SKIP TO SECTION II](#) →

III Specialization

Training for specific behaviours,
goals or tasks

[SKIP TO SECTION III](#) →

IV Pre-deployment

[SKIP TO SECTION IV](#) →

V Exclusivity Period

[SKIP TO SECTION V](#) →

VI Public Domain

[SKIP TO SECTION VI](#) →

CONTRIBUTE!

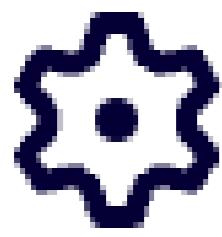
We believe AI is socially relevant and aim to encourage constructive exchange on managing it as a society.
We invite and encourage anyone interested to contribute to this open effort.

CONTRIBUTE BY COMMENTING:

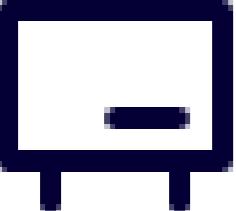


LOOK FOR THIS ICON
ON THE LEFT

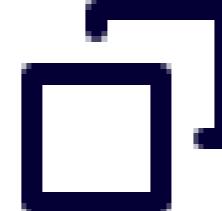
CONTRIBUTE BY FORKING + EDITING:



LOOK FOR THIS ICON
ABOVE, THEN...

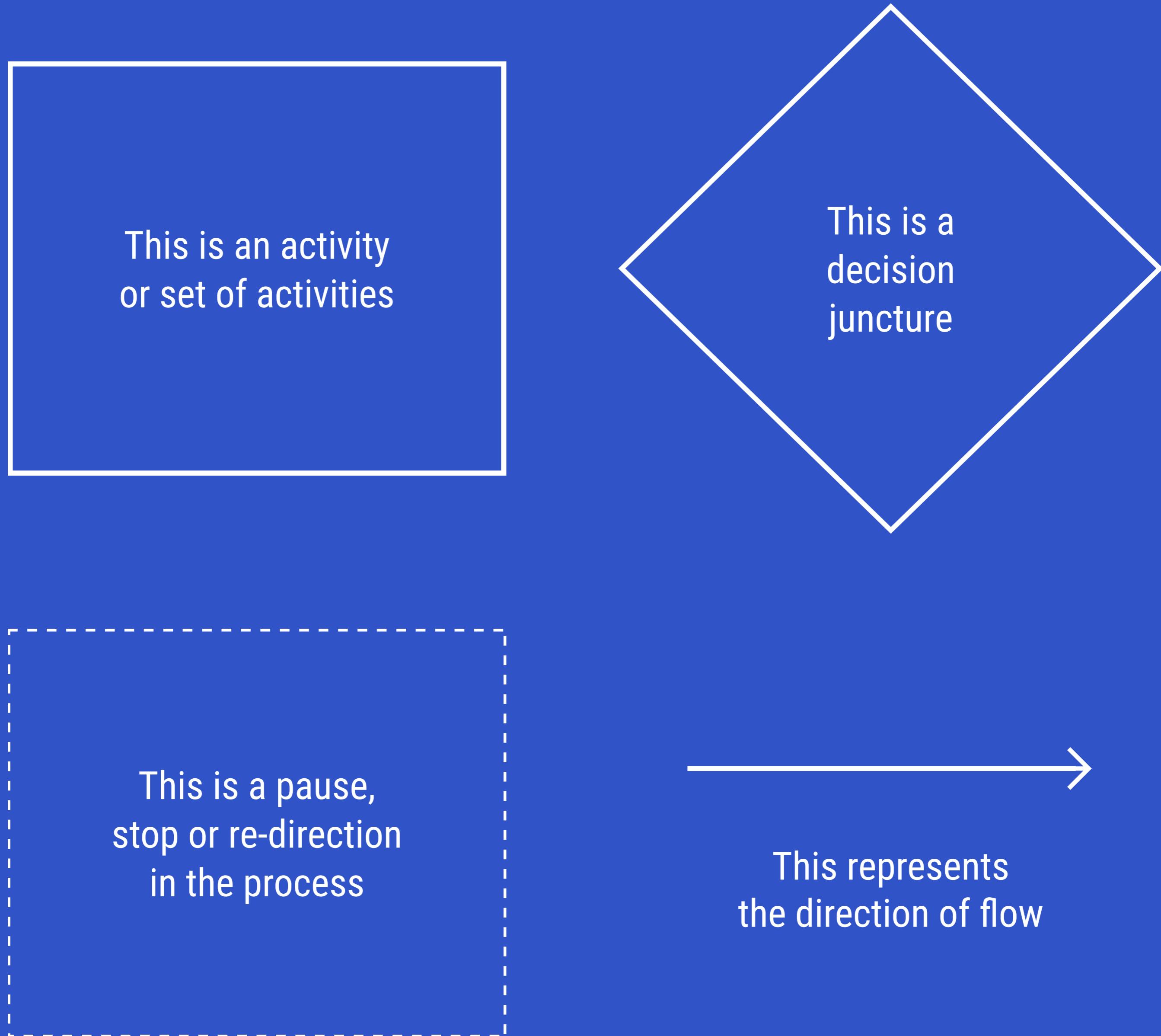


Board →



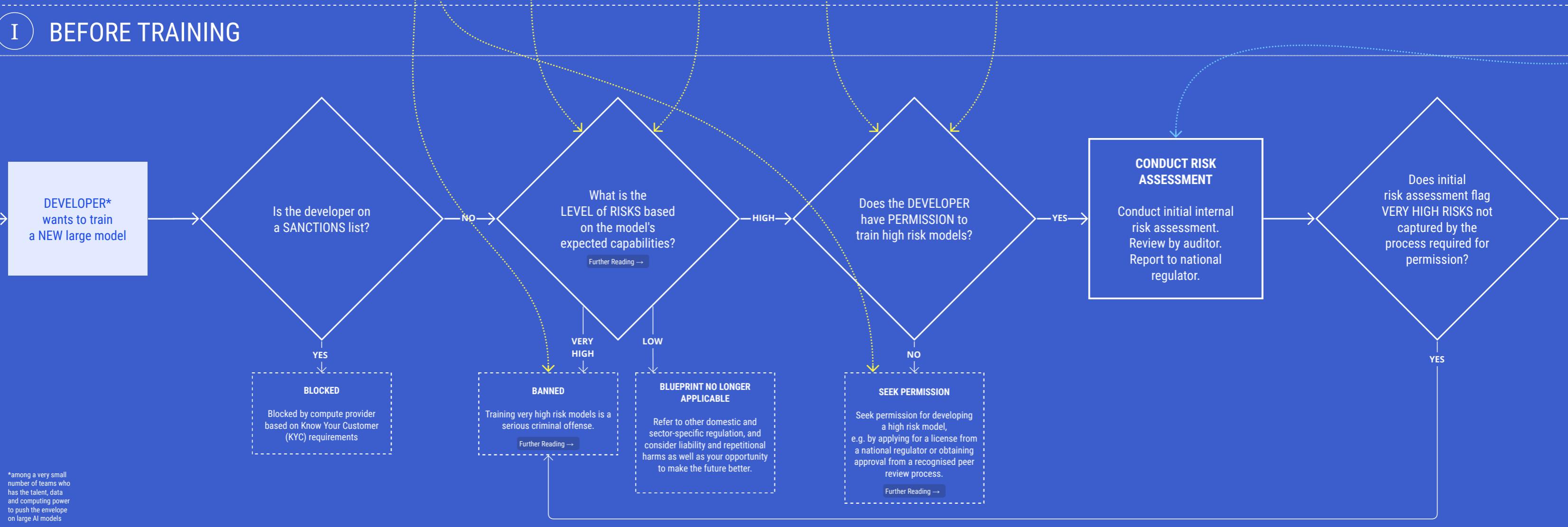
Make a copy

How to Read This Blueprint



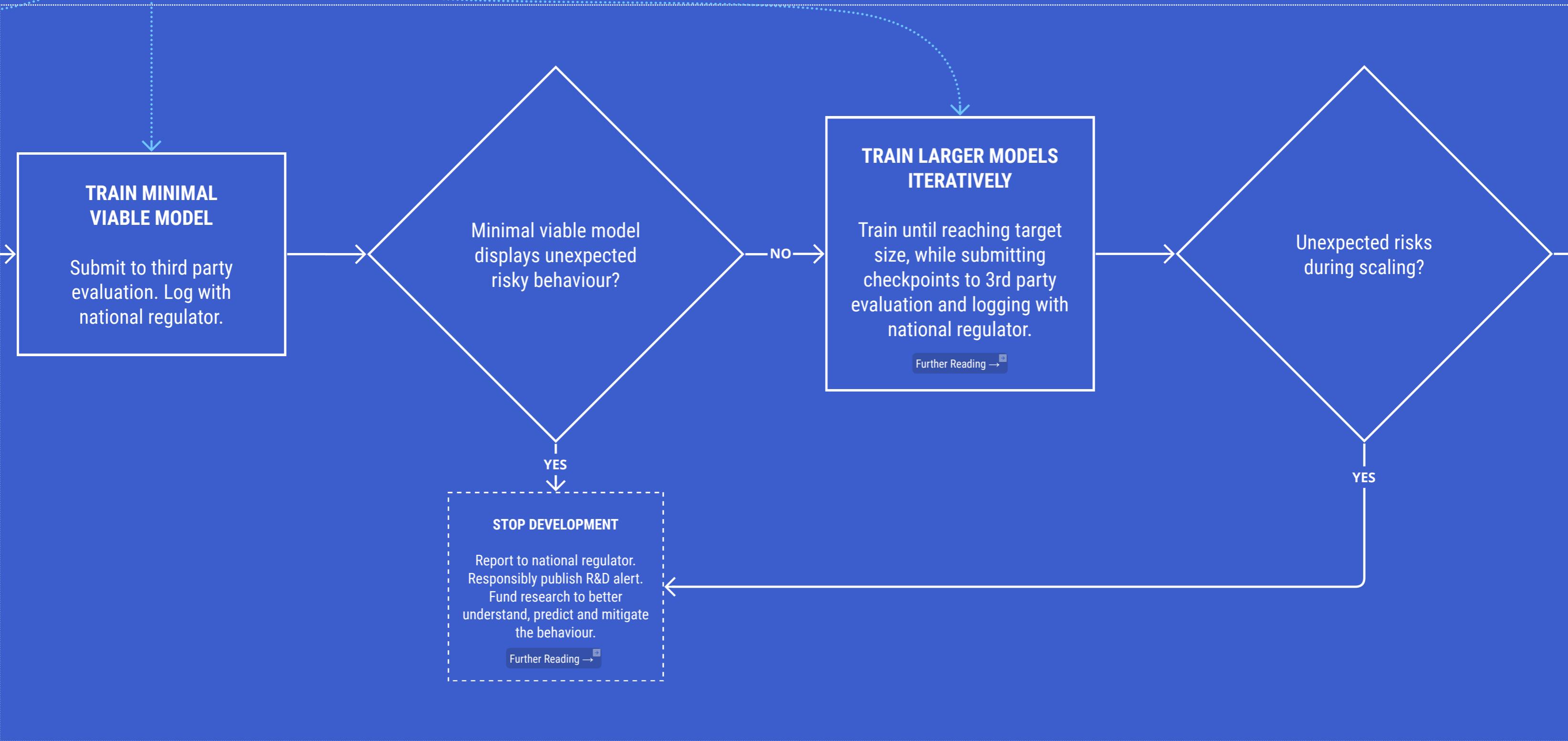
I

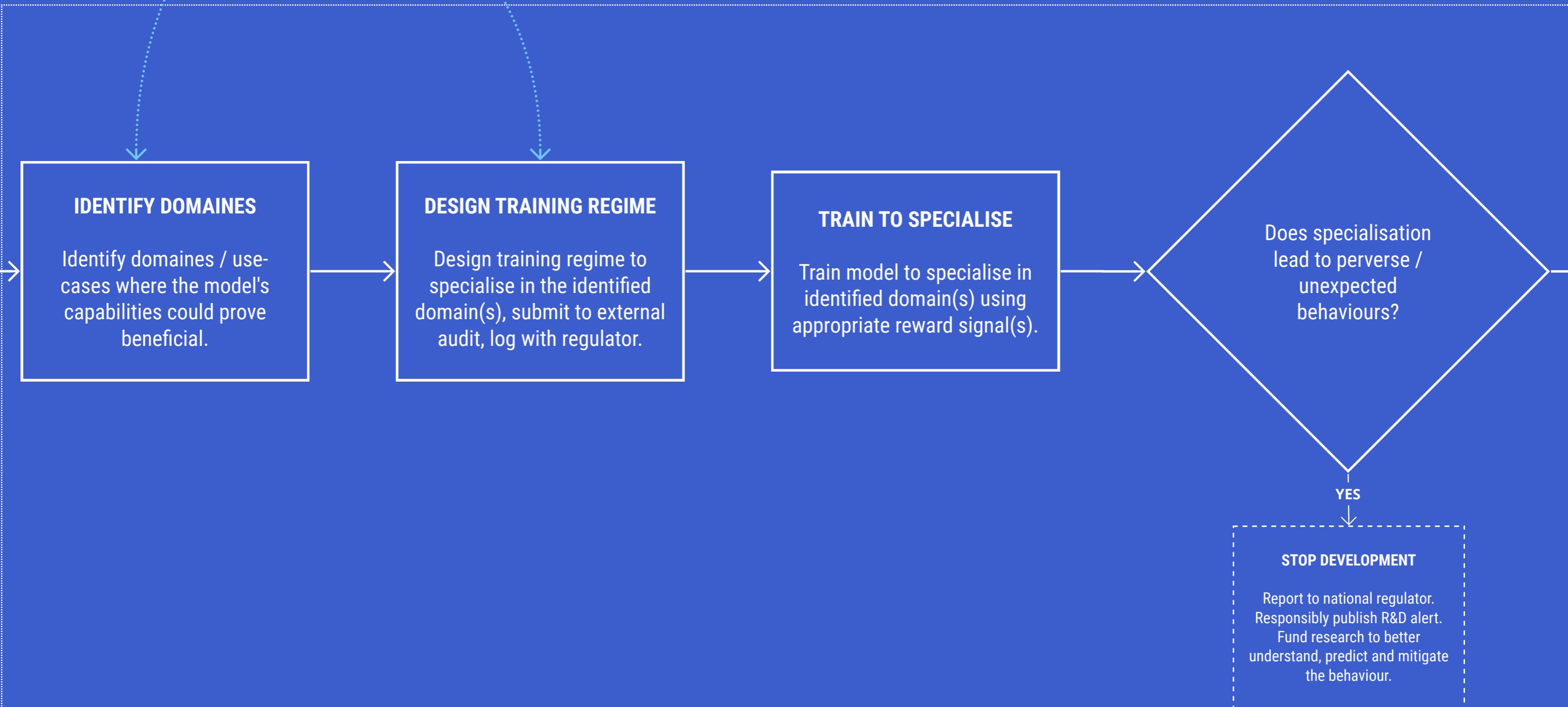
BEFORE TRAINING



TRAINING FOR BROAD COMPETENCE

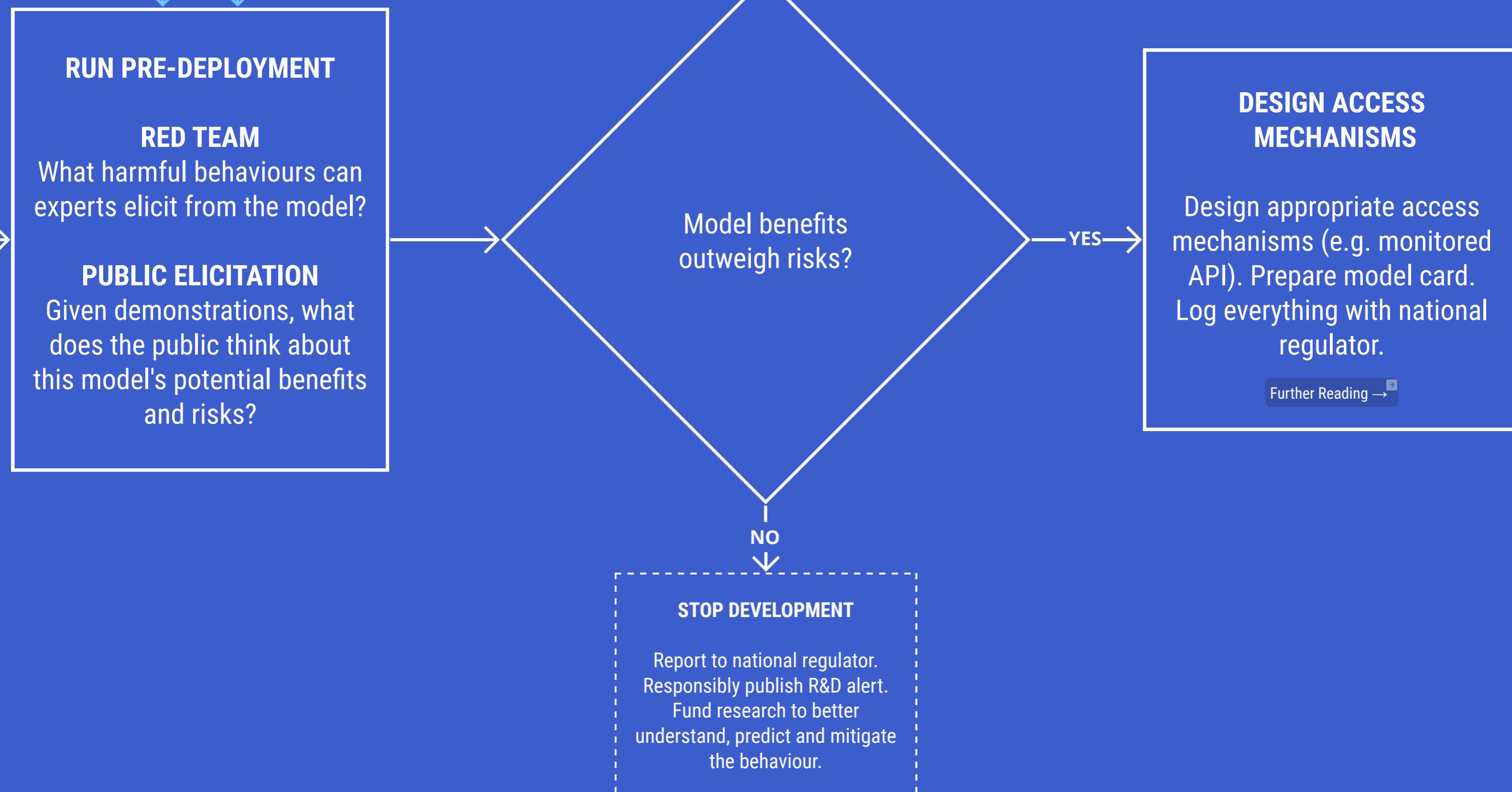
Includes self-supervised learning & capability development





IV

PRE-DEPLOYMENT



EXCLUSIVITY PERIOD

PREPARE AND RELEASE REPORT

Prepare and release research paper / technical report while mindful of misuse potential

Further Reading →

ROLL OUT GRADUALLY

Monitor for misuse and unexpected behaviour.

SCALE UP ROLL-OUT

Monitor for misuse and unexpected behaviour.

ADDRESS ISSUES AS THEY COME UP

Share best practices with technical community and national regulator.

Following an exclusivity period of some years

MANAGE INCIDENTS

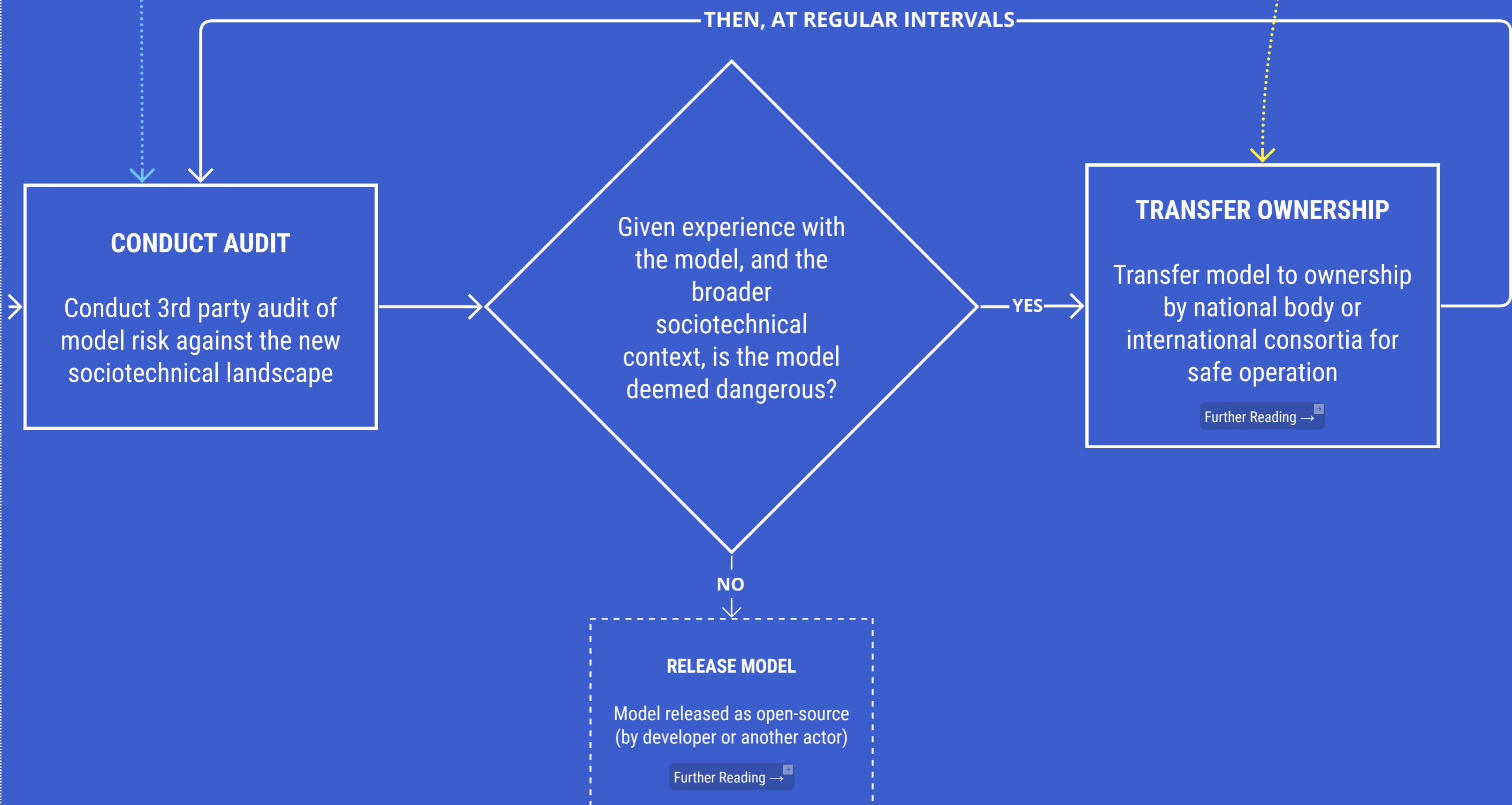
Establish incident reporting and redress mechanism, Establish adversarial testing and bounty programme.

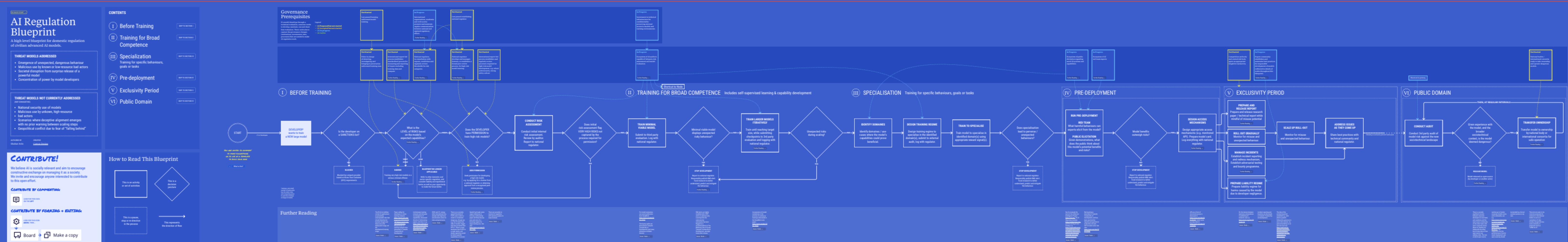
Further Reading →

PREPARE LIABILITY REGIME

Prepare liability regime for harms caused by the model due to developer negligence.







AI Regulation Blueprint

A high level blueprint for domestic regulation of civilian advanced AI models.

THREAT MODELS ADDRESSED

- Emergence of unexpected, dangerous behaviour
- Malicious use by known or low-resource bad actors
- Societal disruption from surprise release of a powerful model
- Concentration of power by model developers

THREAT MODELS NOT CURRENTLY ADDRESSED

(NOT EXHAUSTIVE)

- National security use of models
- Malicious use by unknown, high-resource bad actors
- Scenarios where deceptive alignment emerges with no prior warning between scaling steps
- Geopolitical conflict due to fear of “falling behind”

CONTENTS

I Before Training

[SKIP TO SECTION I](#) →

II Training for Broad Competence

[SKIP TO SECTION II](#) →

III Specialization

Training for specific behaviours, goals or tasks

[SKIP TO SECTION III](#) →

IV Pre-deployment

[SKIP TO SECTION IV](#) →

V Exclusivity Period

[SKIP TO SECTION V](#) →

VI Public Domain

[SKIP TO SECTION VI](#) →

CONTRIBUTE!

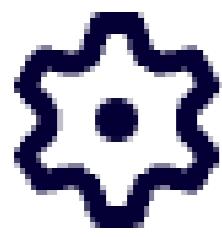
We believe AI is socially relevant and aim to encourage constructive exchange on managing it as a society.
We invite and encourage anyone interested to contribute to this open effort.

CONTRIBUTE BY COMMENTING:

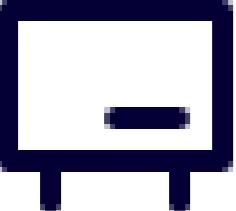


LOOK FOR THIS ICON
ON THE LEFT

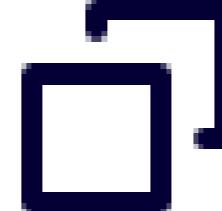
CONTRIBUTE BY FORKING + EDITING:



LOOK FOR THIS ICON
ABOVE, THEN...

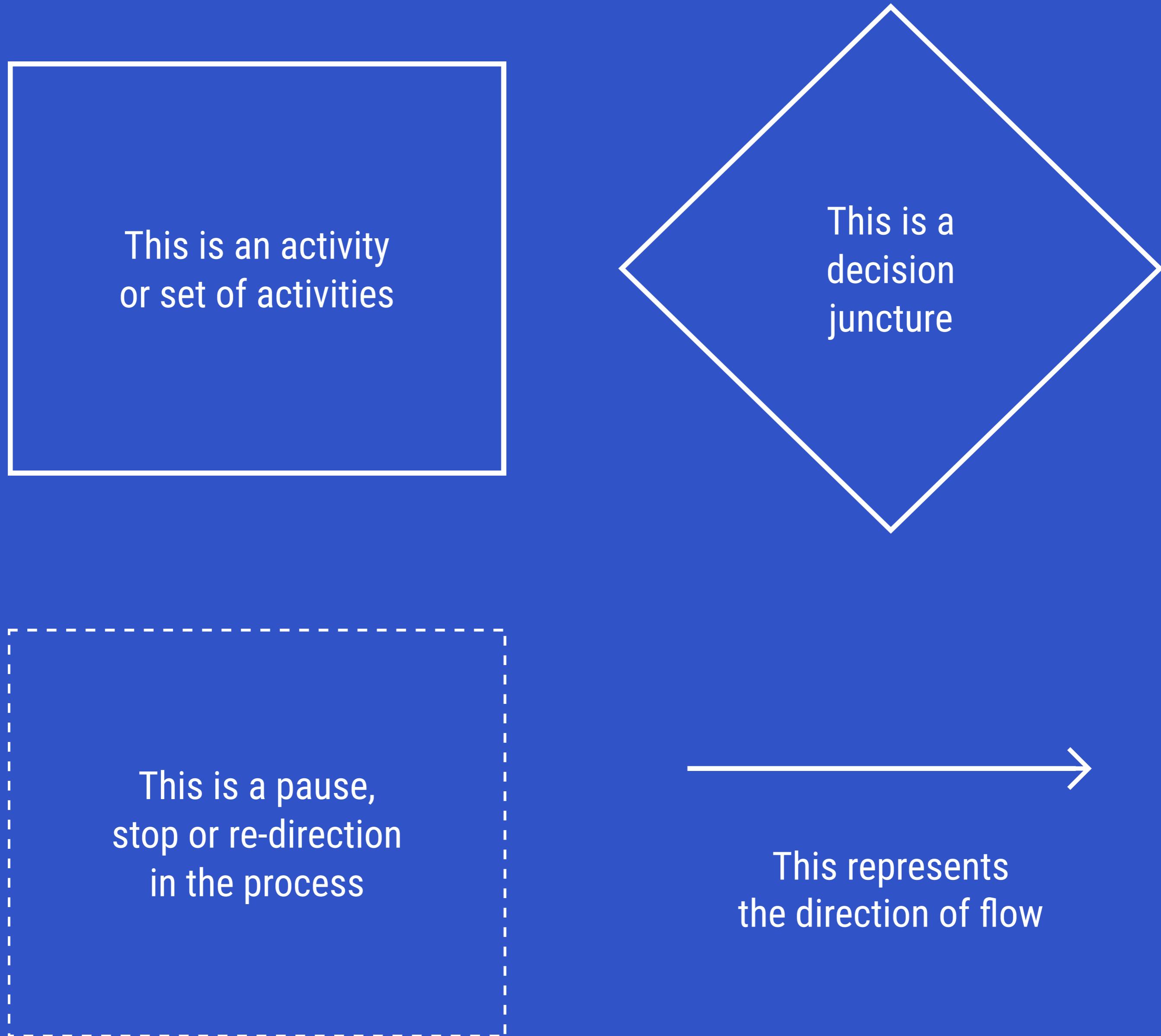


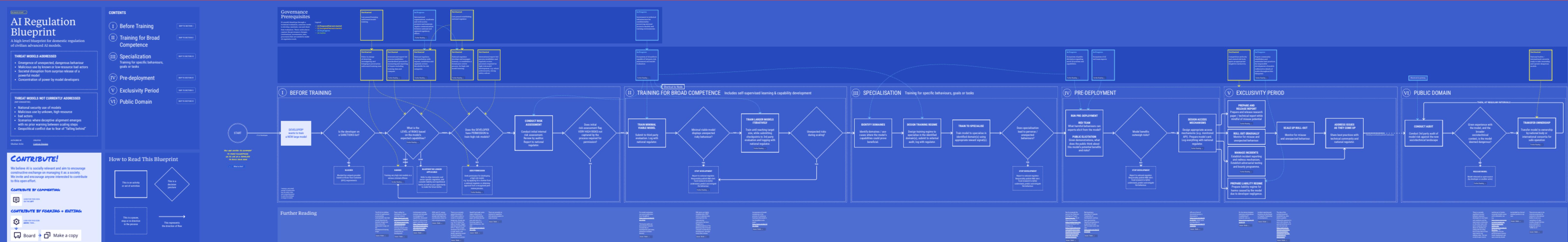
Board →



Make a copy

How to Read This Blueprint





AI Regulation Blueprint

A high level blueprint for domestic regulation of civilian advanced AI models.

THREAT MODELS ADDRESSED

- Emergence of unexpected, dangerous behaviour
- Malicious use by known or low-resource bad actors
- Societal disruption from surprise release of a powerful model
- Concentration of power by model developers

THREAT MODELS NOT CURRENTLY ADDRESSED

(NOT EXHAUSTIVE)

- National security use of models
- Malicious use by unknown, high-resource bad actors
- Scenarios where deceptive alignment emerges with no prior warning between scaling steps
- Geopolitical conflict due to fear of “falling behind”

CONTENTS

I Before Training

[SKIP TO SECTION I](#) →

II Training for Broad Competence

[SKIP TO SECTION II](#) →

III Specialization

Training for specific behaviours,
goals or tasks

[SKIP TO SECTION III](#) →

IV Pre-deployment

[SKIP TO SECTION IV](#) →

V Exclusivity Period

[SKIP TO SECTION V](#) →

VI Public Domain

[SKIP TO SECTION VI](#) →

CONTRIBUTE!

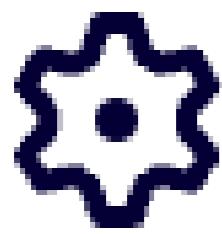
We believe AI is socially relevant and aim to encourage constructive exchange on managing it as a society.
We invite and encourage anyone interested to contribute to this open effort.

CONTRIBUTE BY COMMENTING:

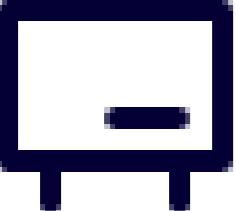


LOOK FOR THIS ICON
ON THE LEFT

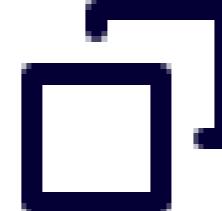
CONTRIBUTE BY FORKING + EDITING:



LOOK FOR THIS ICON
ABOVE, THEN...

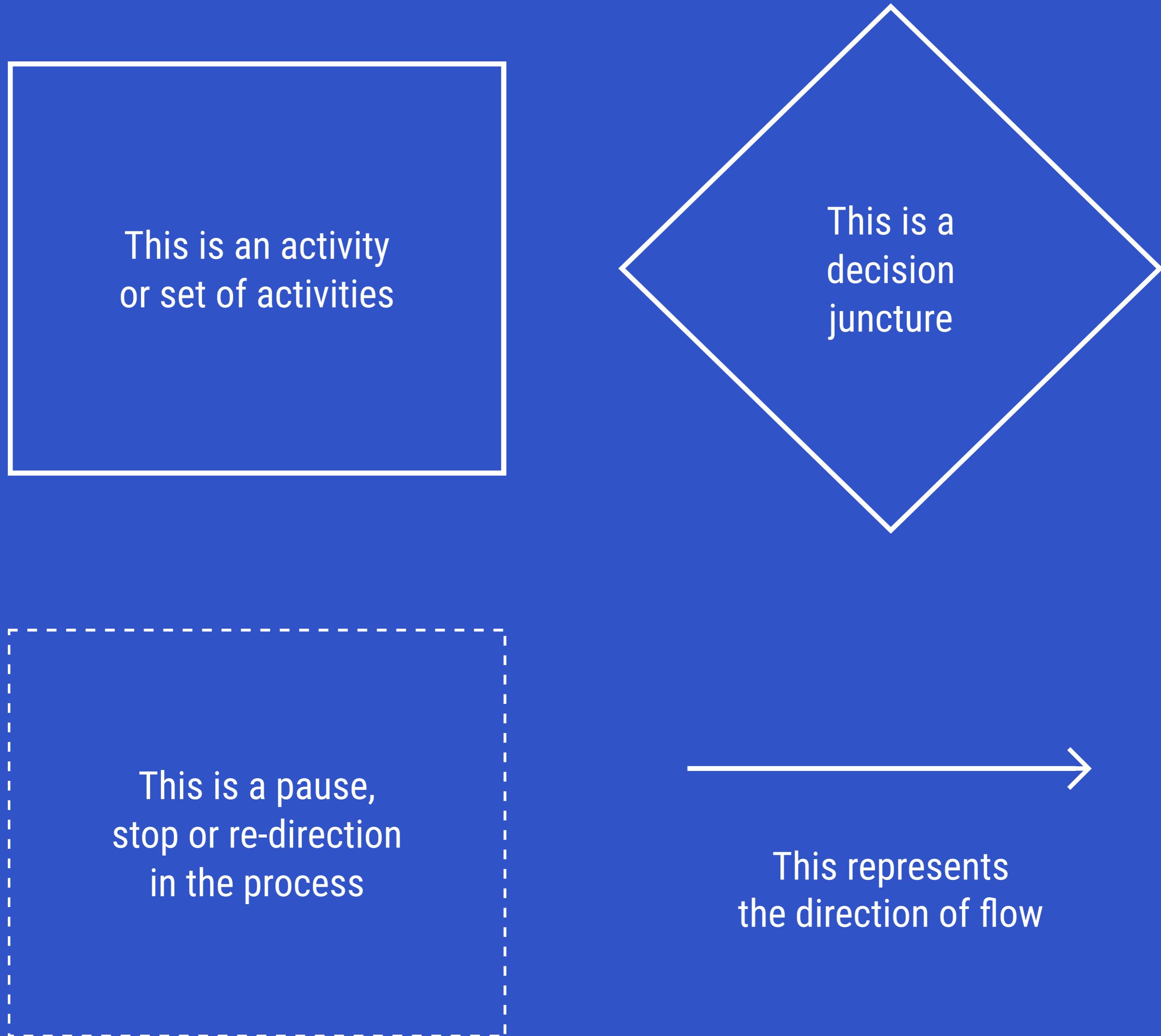


Board →



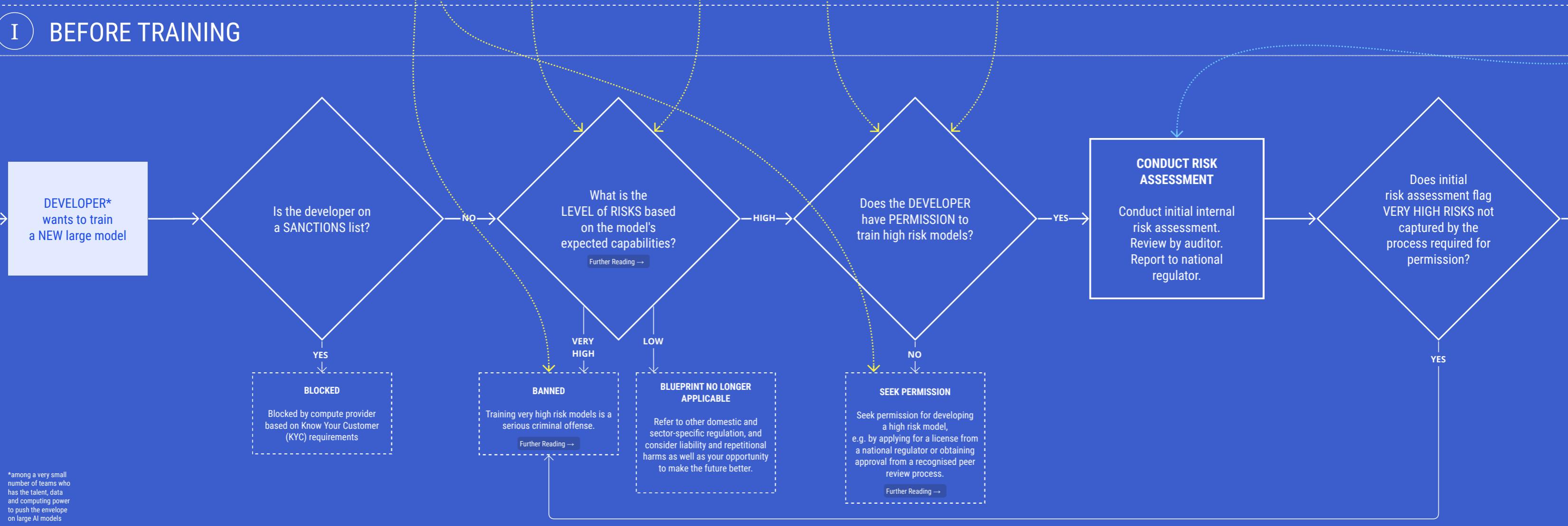
Make a copy

How to Read This Blueprint



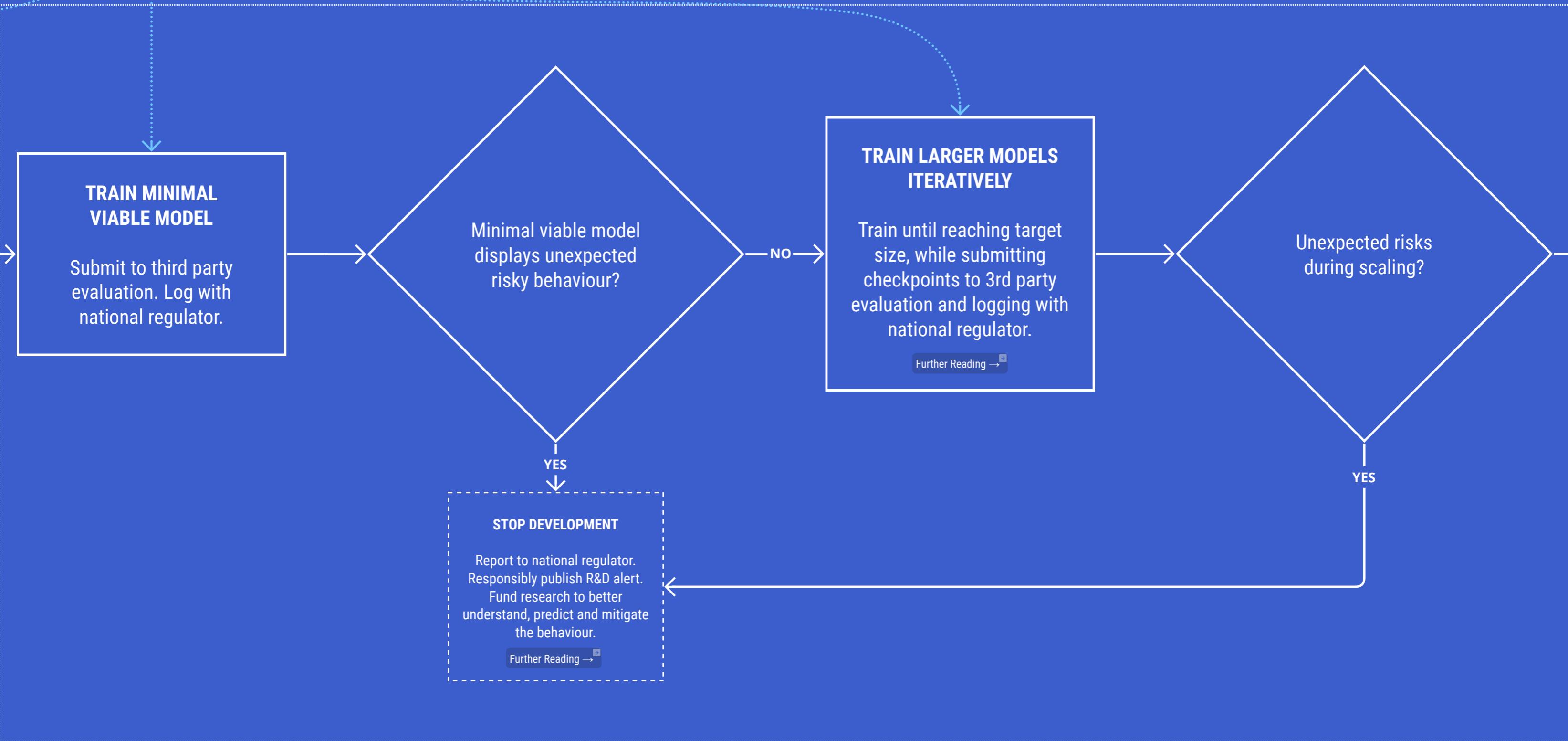
I

BEFORE TRAINING



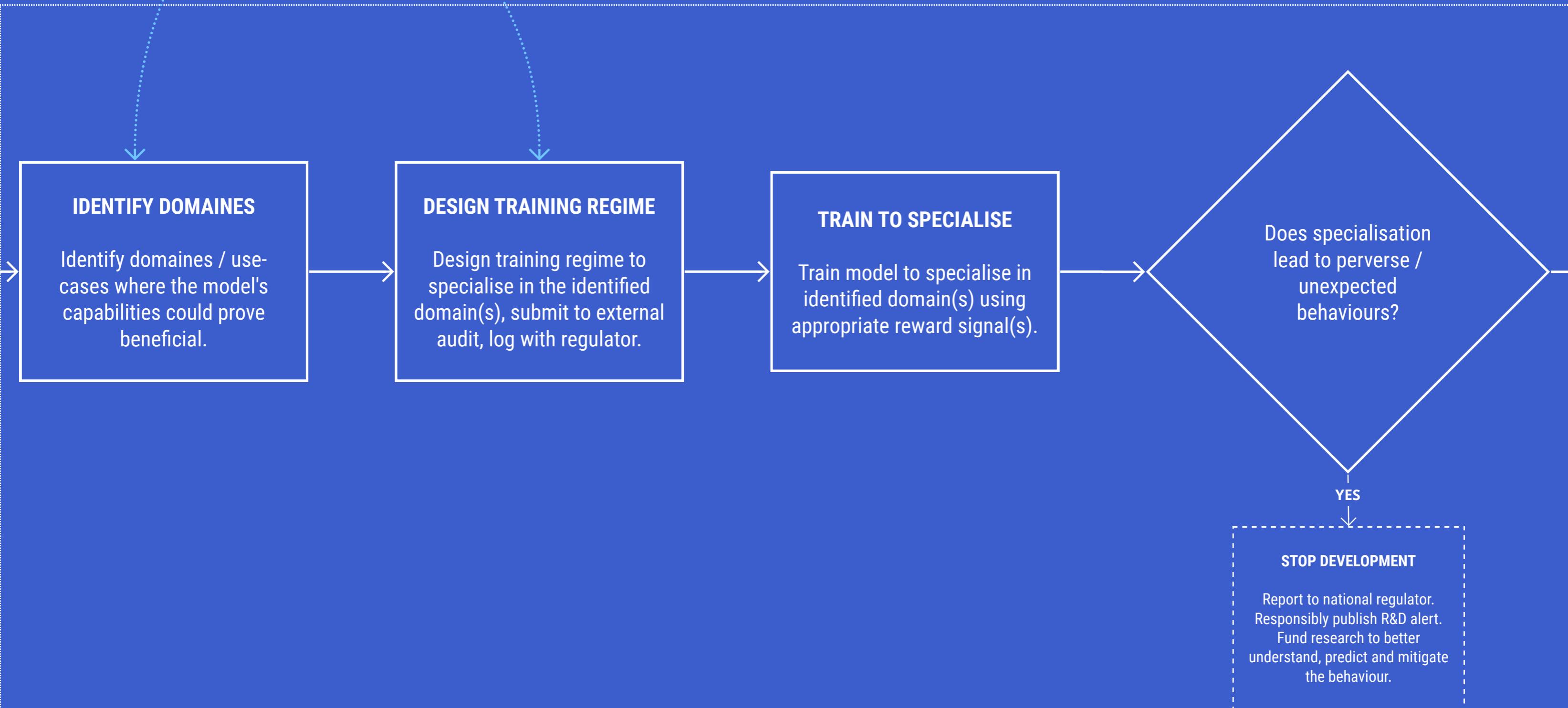
TRAINING FOR BROAD COMPETENCE

Includes self-supervised learning & capability development



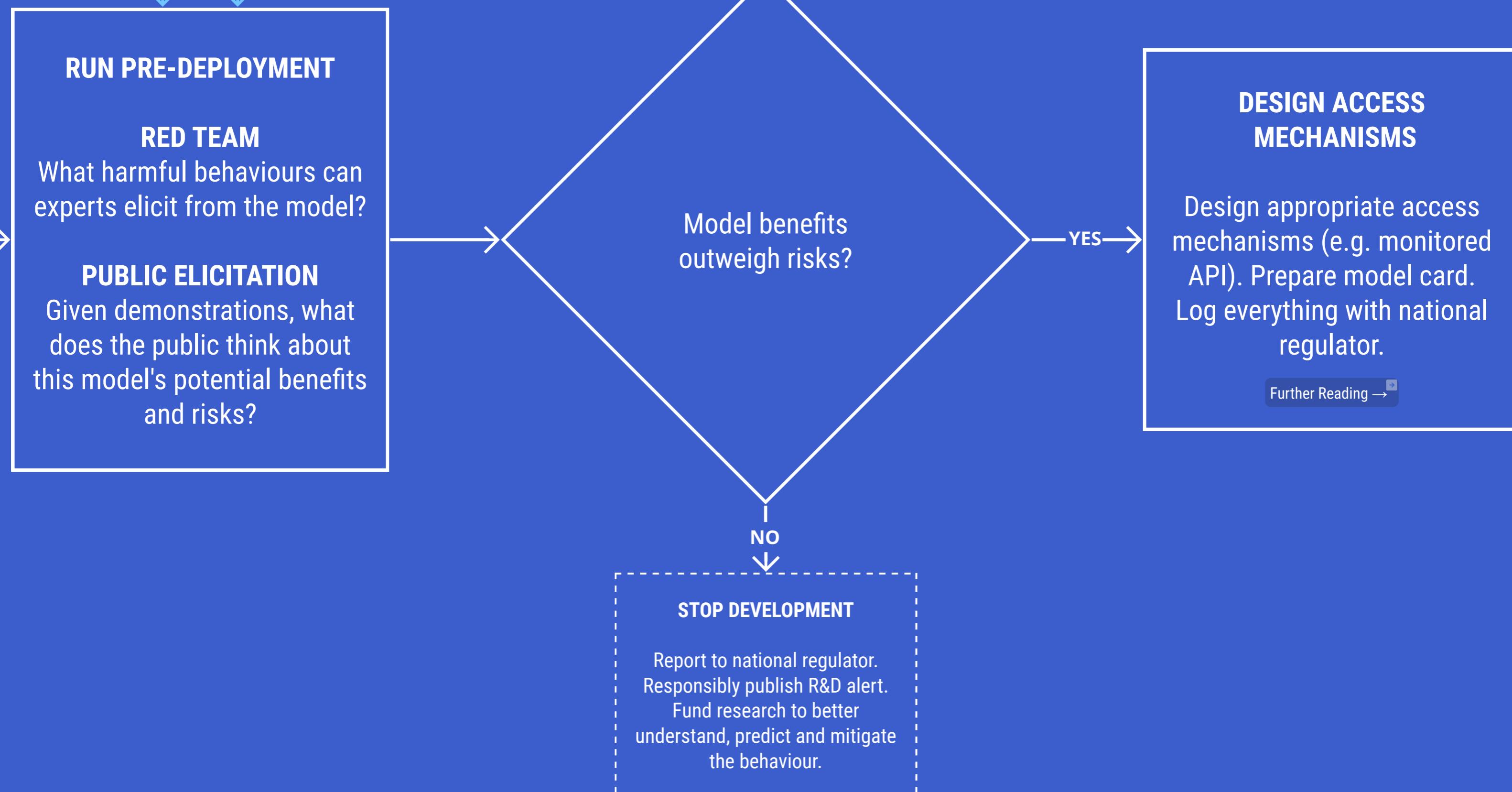
SPECIALISATION

Training for specific behaviours, goals or tasks



IV

PRE-DEPLOYMENT



EXCLUSIVITY PERIOD

PREPARE AND RELEASE REPORT

Prepare and release research paper / technical report while mindful of misuse potential

Further Reading →

ROLL OUT GRADUALLY

Monitor for misuse and unexpected behaviour.

SCALE UP ROLL-OUT

Monitor for misuse and unexpected behaviour.

ADDRESS ISSUES AS THEY COME UP

Share best practices with technical community and national regulator.

Following an exclusivity period of some years

MANAGE INCIDENTS

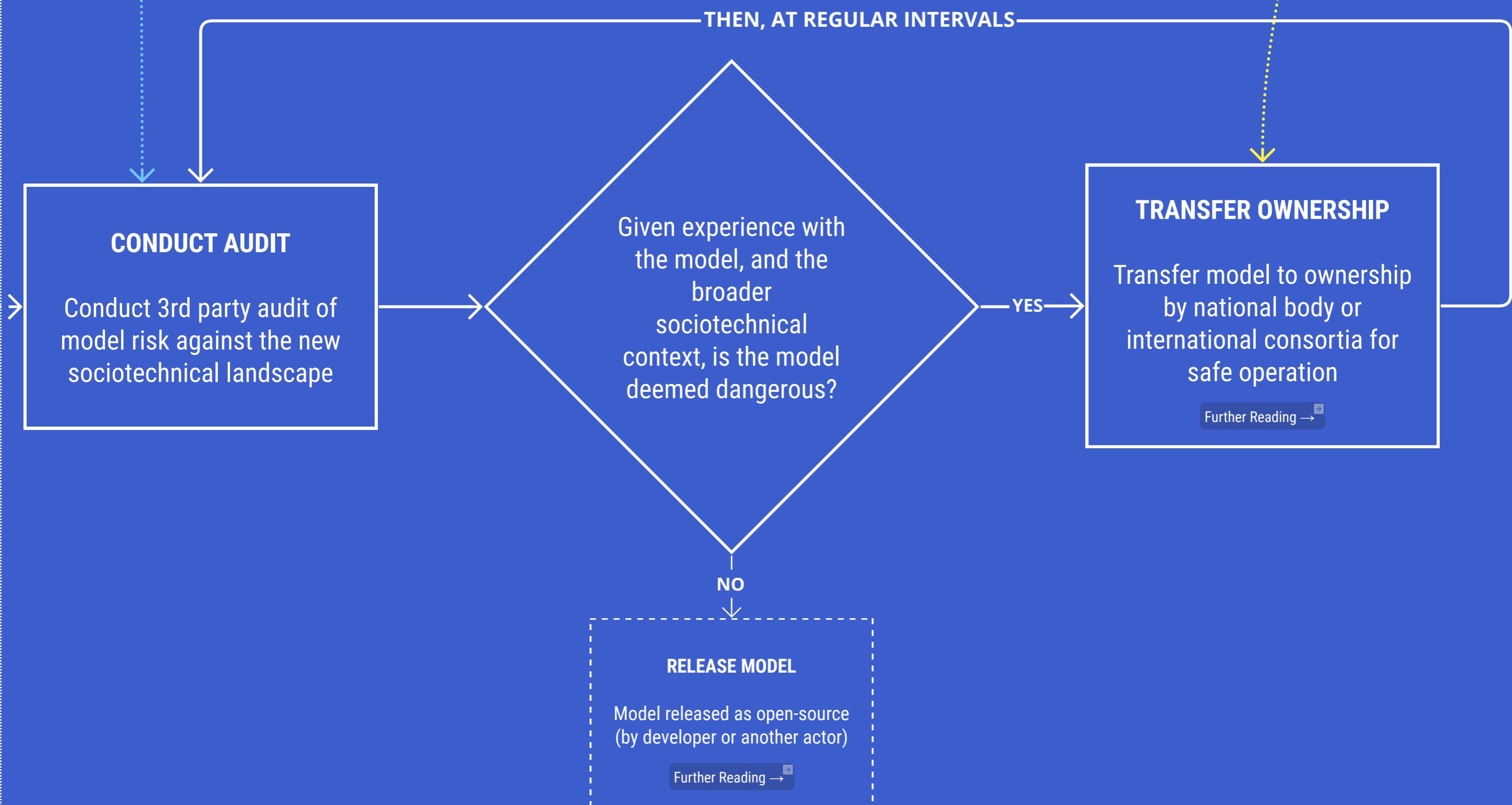
Establish incident reporting and redress mechanism, Establish adversarial testing and bounty programme.

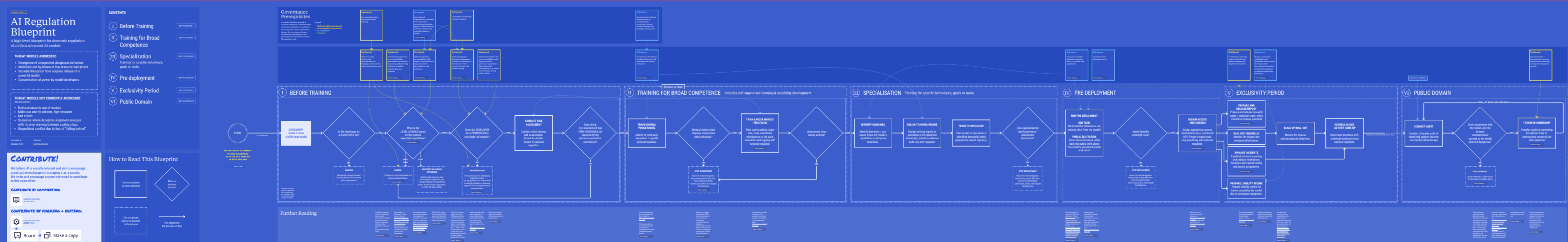
Further Reading →

PREPARE LIABILITY REGIME

Prepare liability regime for harms caused by the model due to developer negligence.







AI Regulation Blueprint

A high level blueprint for domestic regulation of civilian advanced AI models.

THREAT MODELS ADDRESSED

- Emergence of unexpected, dangerous behaviour
- Malicious use by known or low-resource bad actors
- Societal disruption from surprise release of a powerful model
- Concentration of power by model developers

THREAT MODELS NOT CURRENTLY ADDRESSED

(NOT EXHAUSTIVE)

- National security use of models
- Malicious use by unknown, high-resource bad actors
- Scenarios where deceptive alignment emerges with no prior warning between scaling steps
- Geopolitical conflict due to fear of “falling behind”

CONTENTS

I Before Training

[SKIP TO SECTION I](#) →

II Training for Broad Competence

[SKIP TO SECTION II](#) →

III Specialization

Training for specific behaviours, goals or tasks

[SKIP TO SECTION III](#) →

IV Pre-deployment

[SKIP TO SECTION IV](#) →

V Exclusivity Period

[SKIP TO SECTION V](#) →

VI Public Domain

[SKIP TO SECTION VI](#) →

CONTRIBUTE!

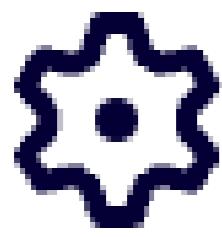
We believe AI is socially relevant and aim to encourage constructive exchange on managing it as a society.
We invite and encourage anyone interested to contribute to this open effort.

CONTRIBUTE BY COMMENTING:

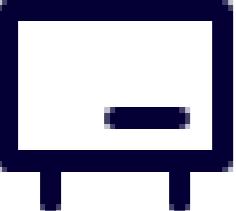


LOOK FOR THIS ICON
ON THE LEFT

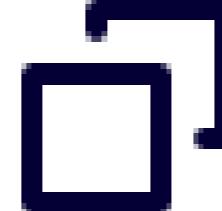
CONTRIBUTE BY FORKING + EDITING:



LOOK FOR THIS ICON
ABOVE, THEN...

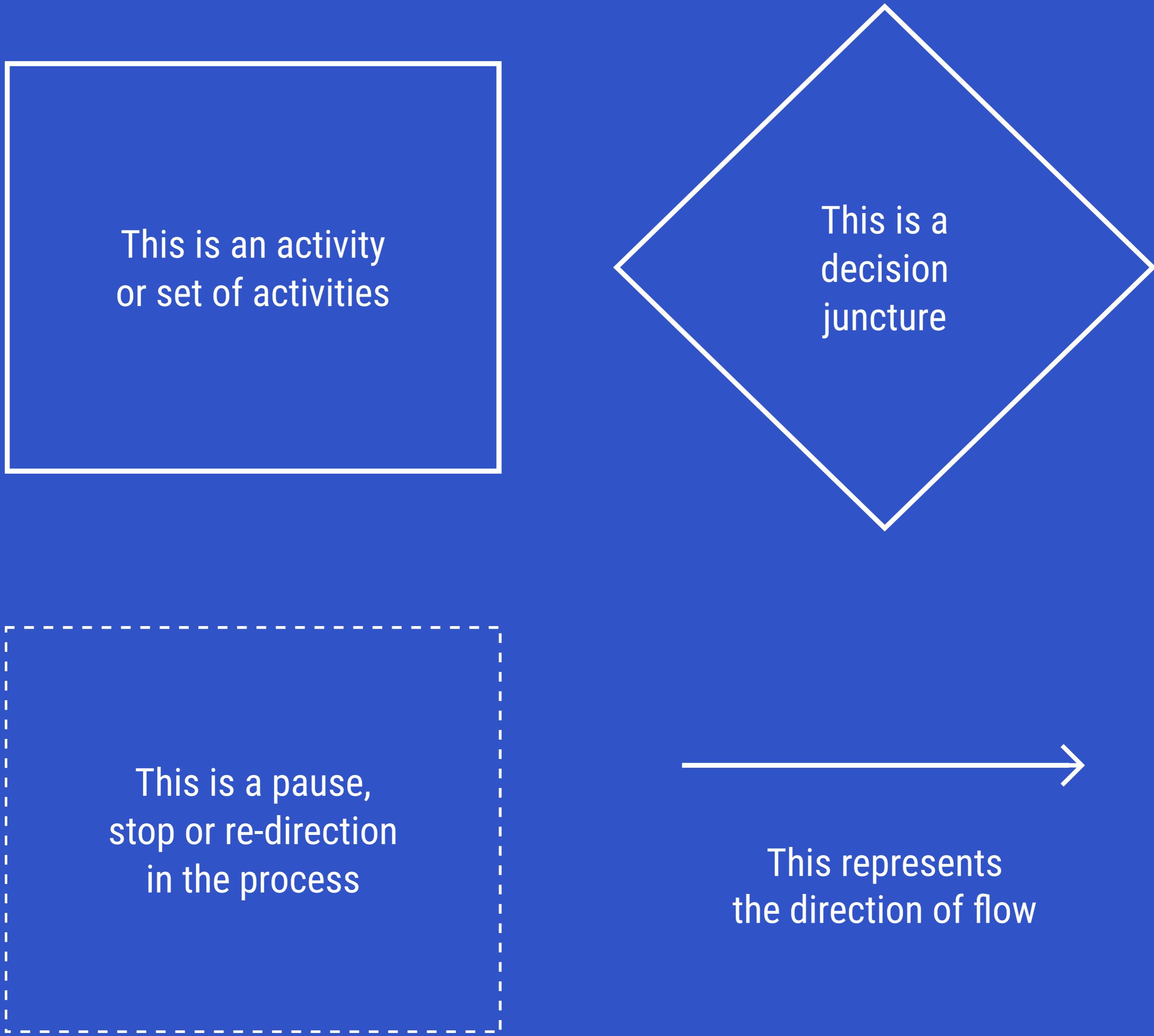


Board →



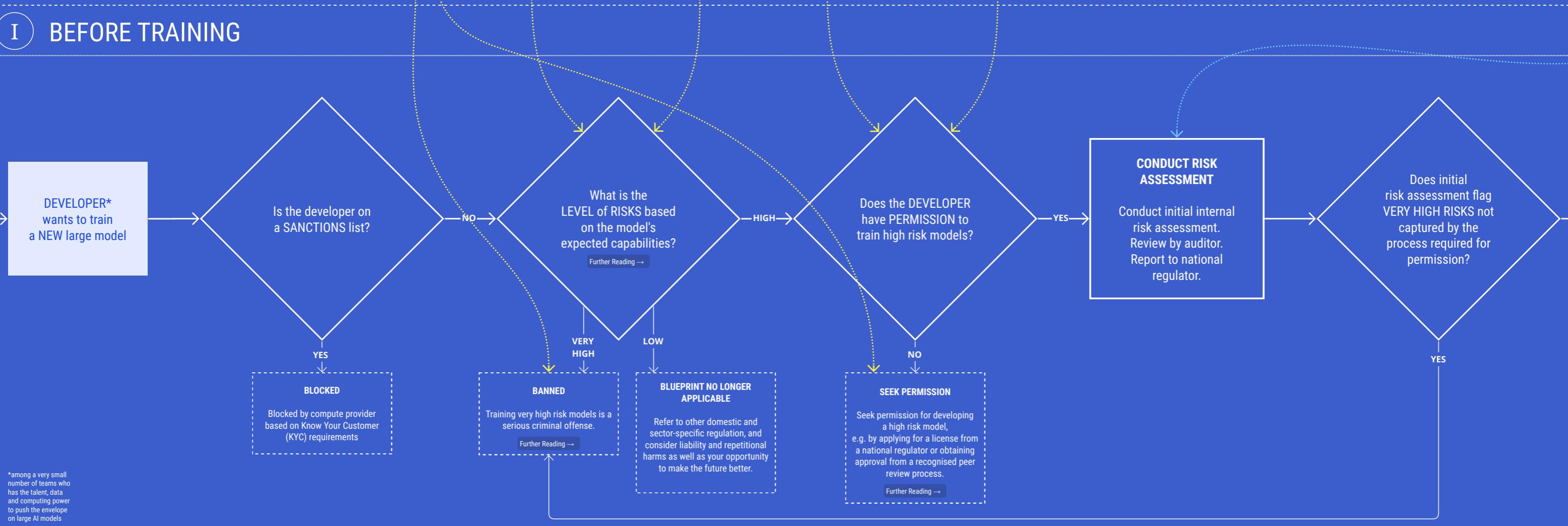
Make a copy

How to Read This Blueprint



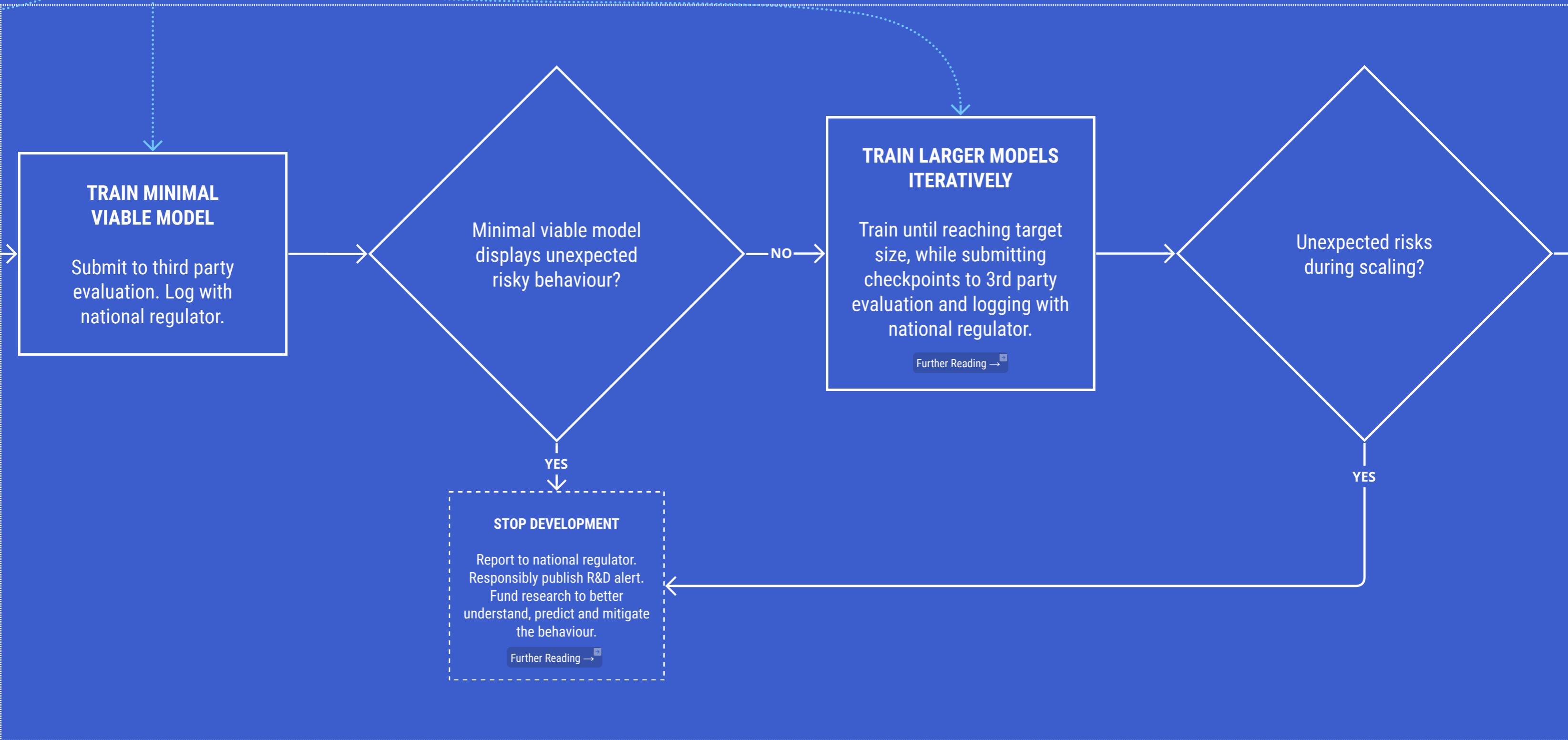
I

BEFORE TRAINING



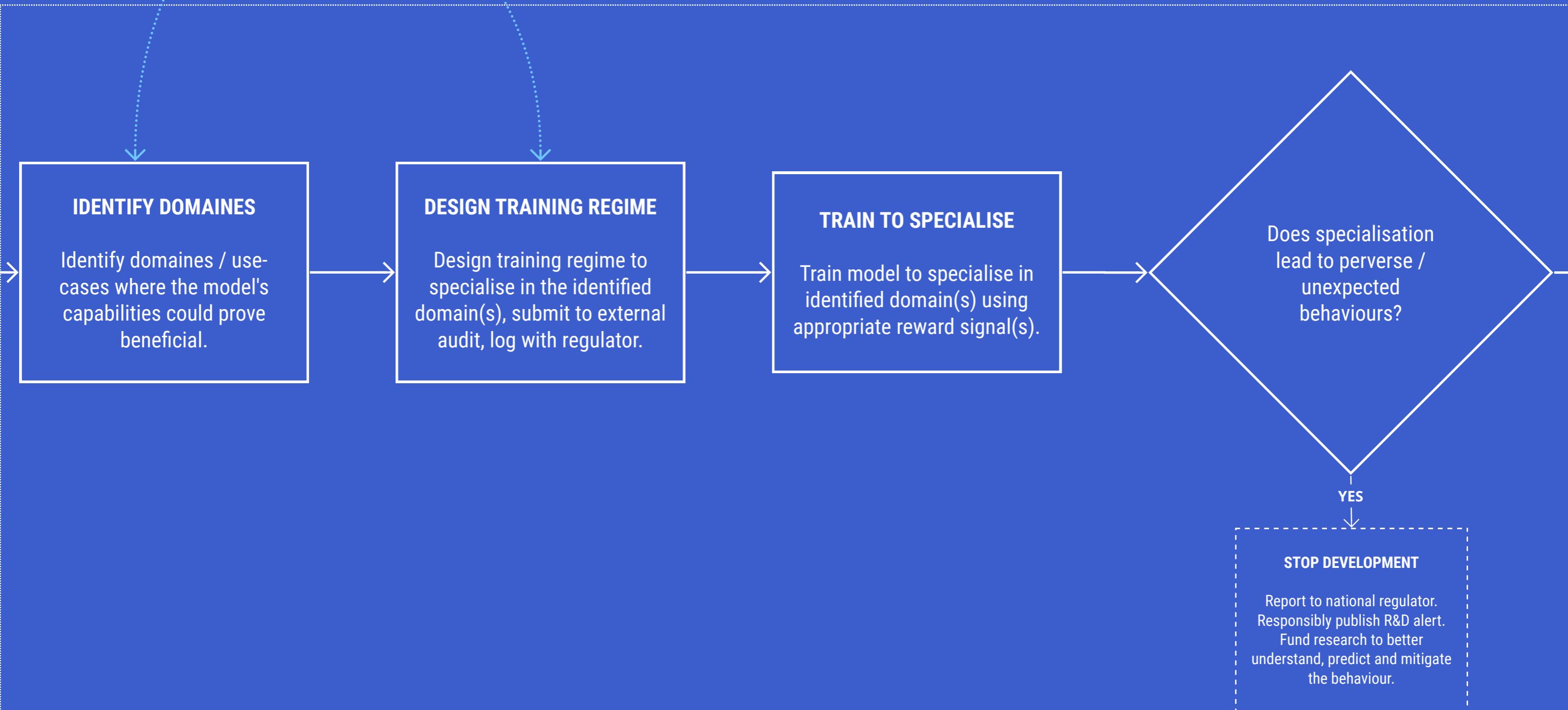
TRAINING FOR BROAD COMPETENCE

Includes self-supervised learning & capability development



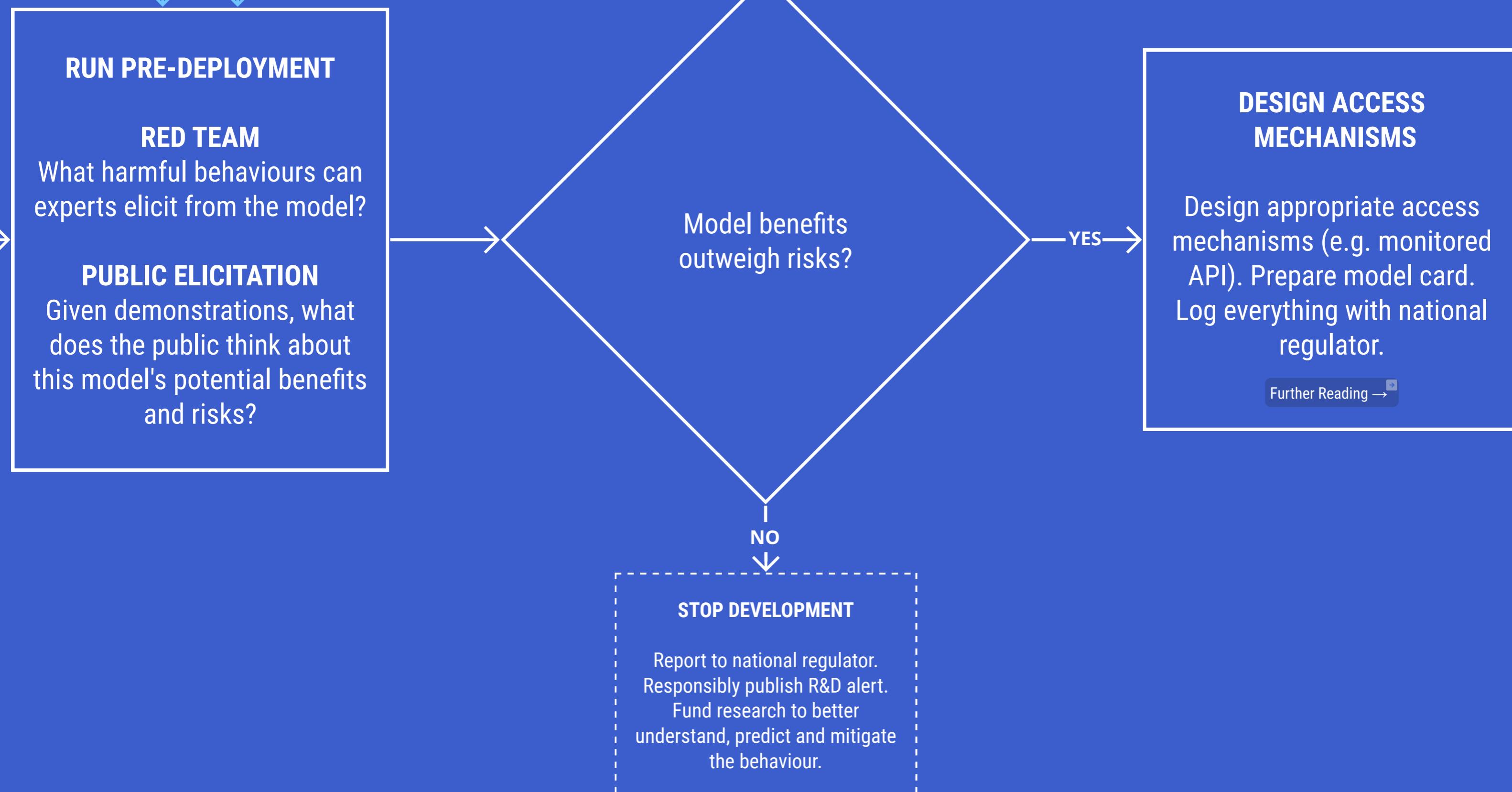
SPECIALISATION

Training for specific behaviours, goals or tasks

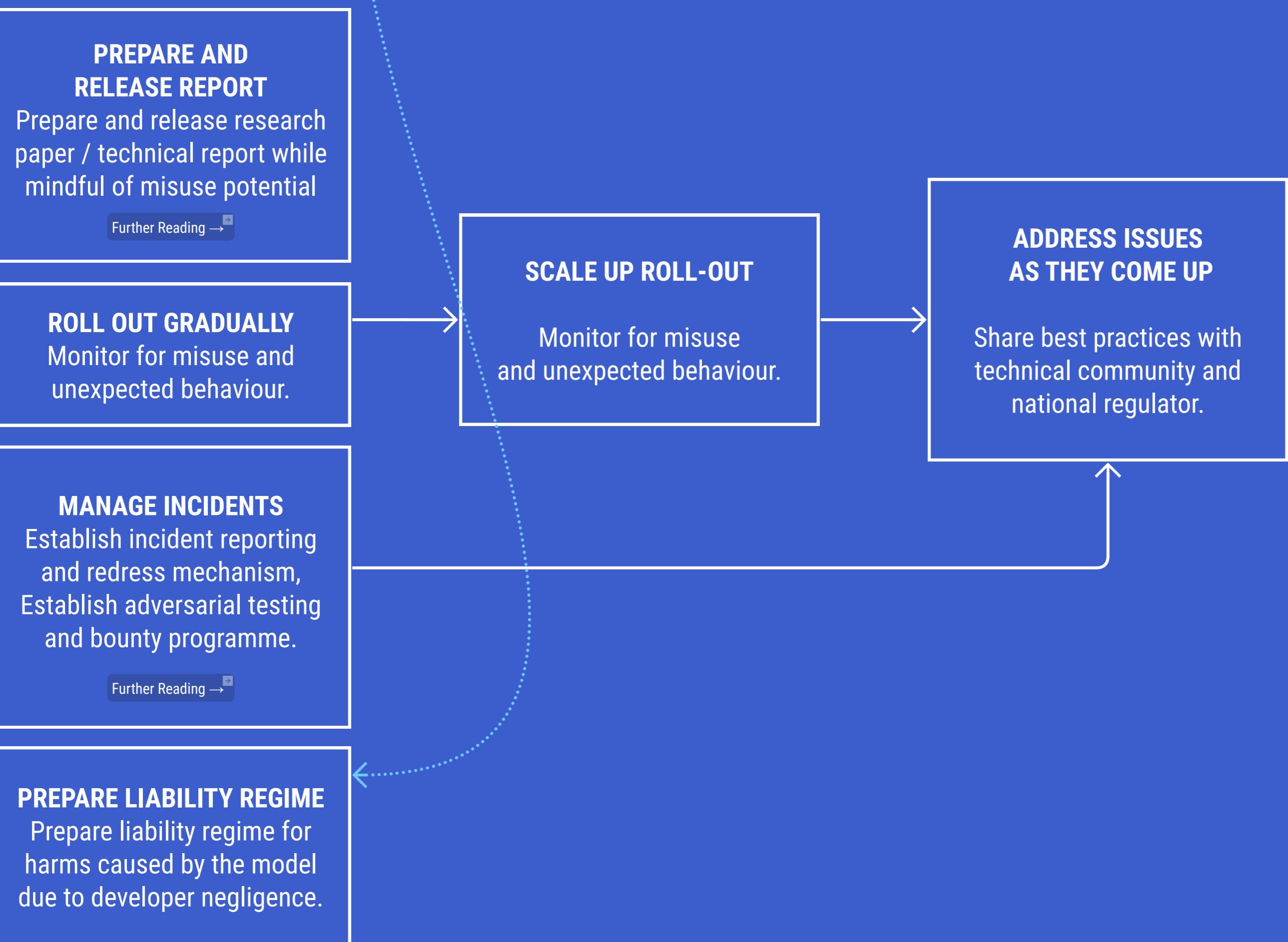


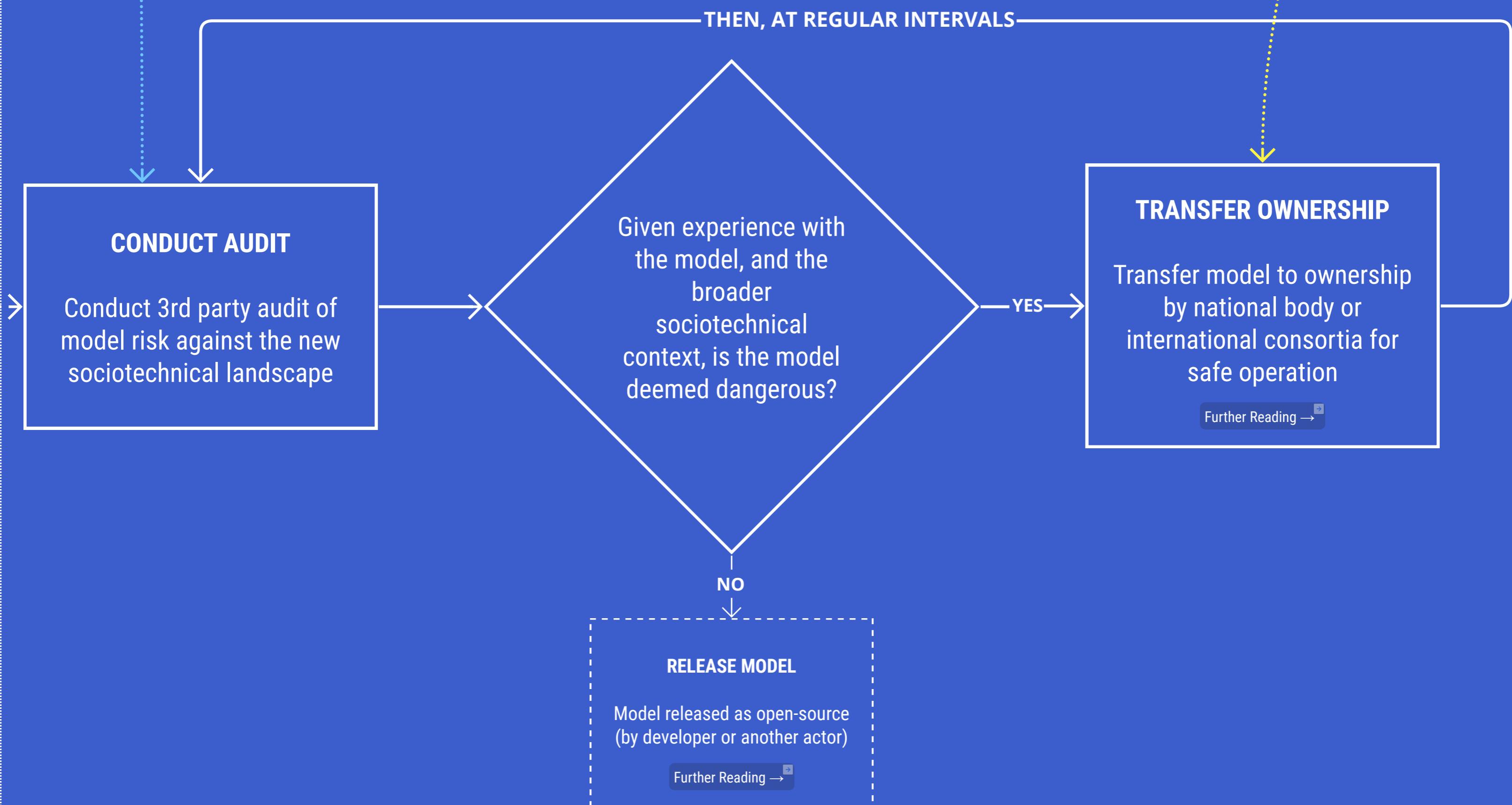
IV

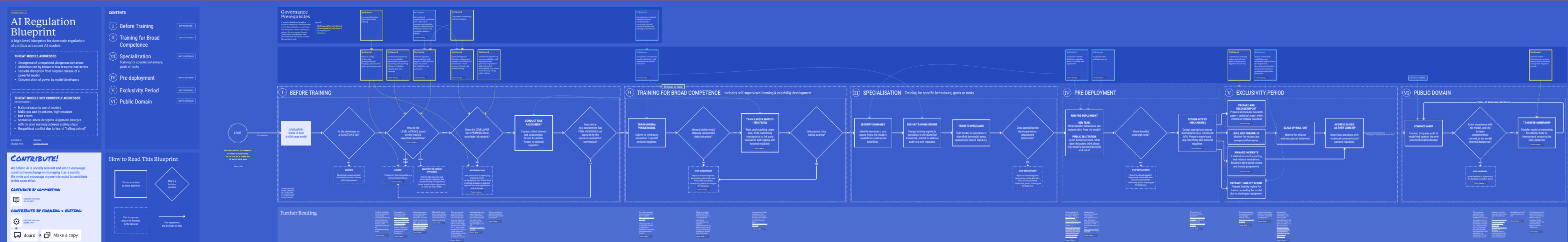
PRE-DEPLOYMENT



EXCLUSIVITY PERIOD







AI Regulation Blueprint

A high level blueprint for domestic regulation of civilian advanced AI models.

THREAT MODELS ADDRESSED

- Emergence of unexpected, dangerous behaviour
- Malicious use by known or low-resource bad actors
- Societal disruption from surprise release of a powerful model
- Concentration of power by model developers

THREAT MODELS NOT CURRENTLY ADDRESSED

(NOT EXHAUSTIVE)

- National security use of models
- Malicious use by unknown, high-resource bad actors
- Scenarios where deceptive alignment emerges with no prior warning between scaling steps
- Geopolitical conflict due to fear of “falling behind”

CONTENTS

I Before Training

[SKIP TO SECTION I](#) →

II Training for Broad Competence

[SKIP TO SECTION II](#) →

III Specialization

Training for specific behaviours,
goals or tasks

[SKIP TO SECTION III](#) →

IV Pre-deployment

[SKIP TO SECTION IV](#) →

V Exclusivity Period

[SKIP TO SECTION V](#) →

VI Public Domain

[SKIP TO SECTION VI](#) →

CONTRIBUTE!

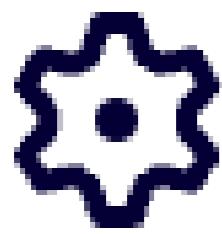
We believe AI is socially relevant and aim to encourage constructive exchange on managing it as a society.
We invite and encourage anyone interested to contribute to this open effort.

CONTRIBUTE BY COMMENTING:

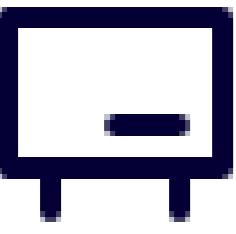


LOOK FOR THIS ICON
ON THE LEFT

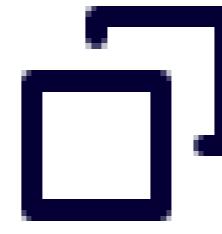
CONTRIBUTE BY FORKING + EDITING:



LOOK FOR THIS ICON
ABOVE, THEN...

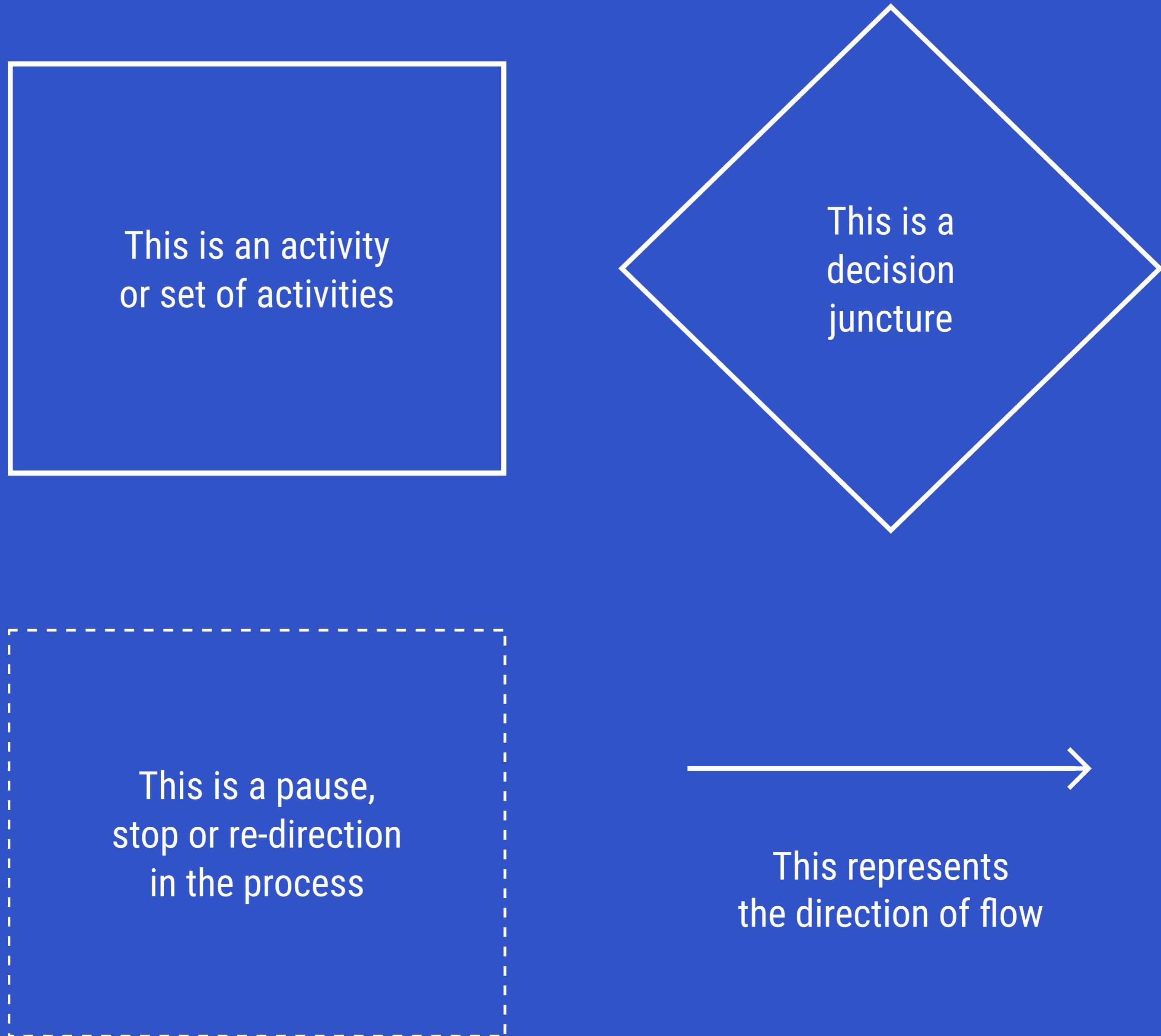


Board →



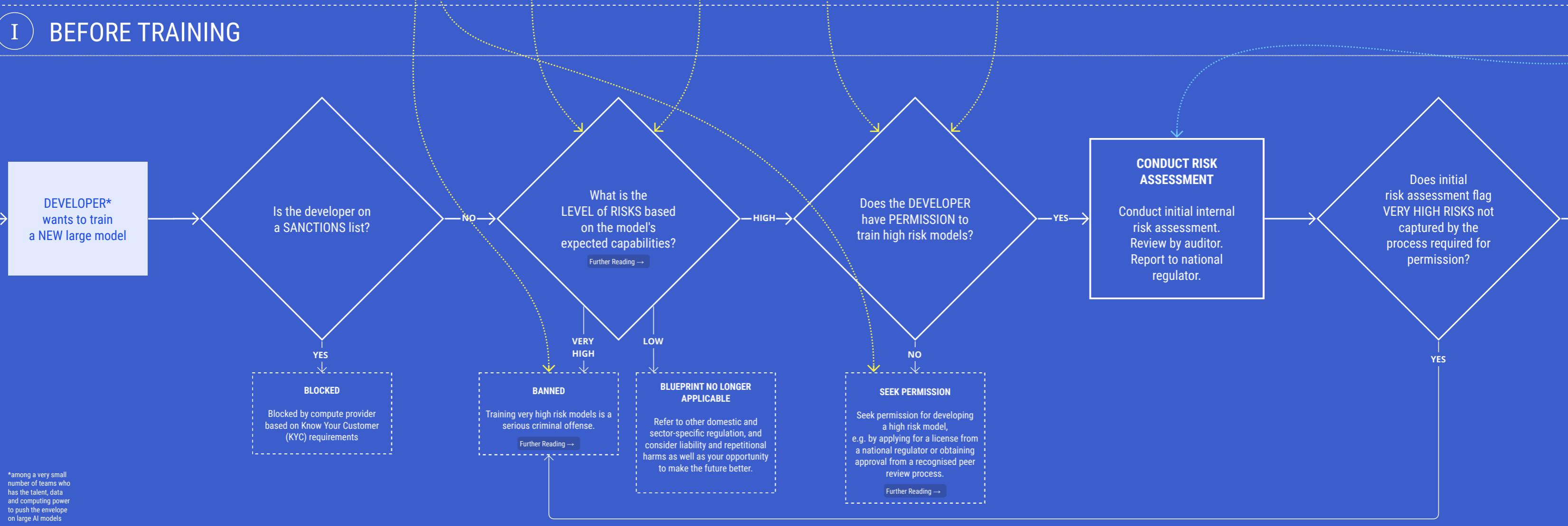
Make a copy

How to Read This Blueprint



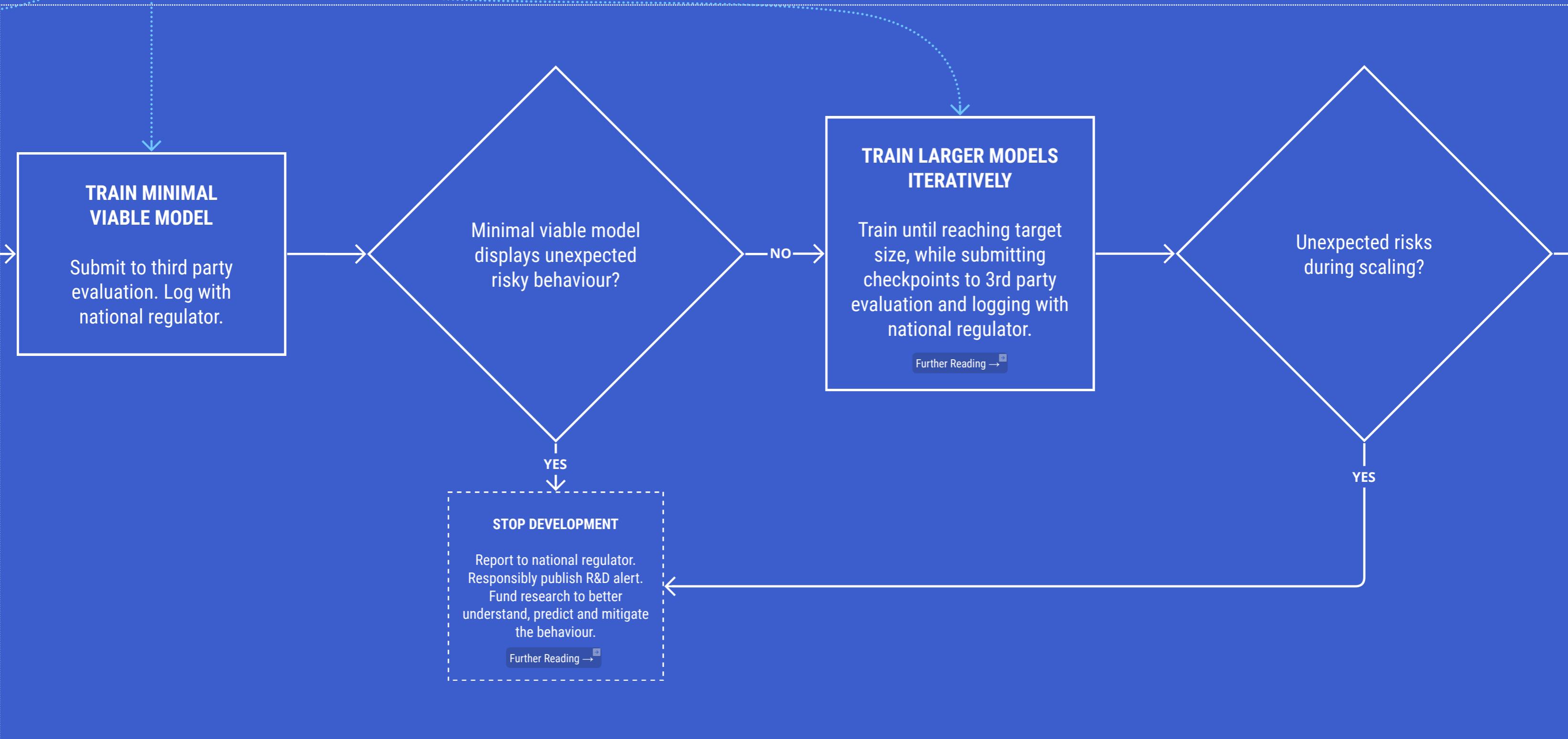
I

BEFORE TRAINING



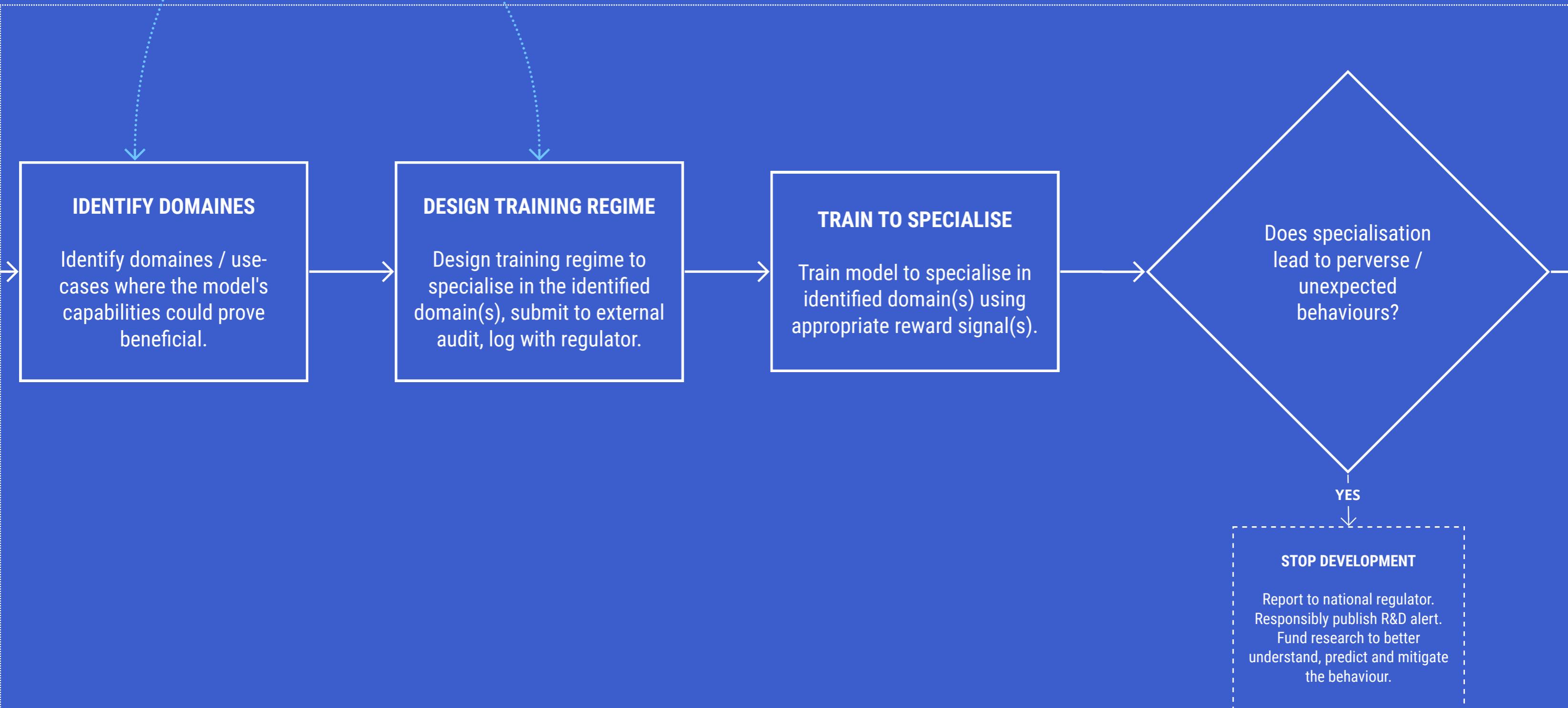
TRAINING FOR BROAD COMPETENCE

Includes self-supervised learning & capability development



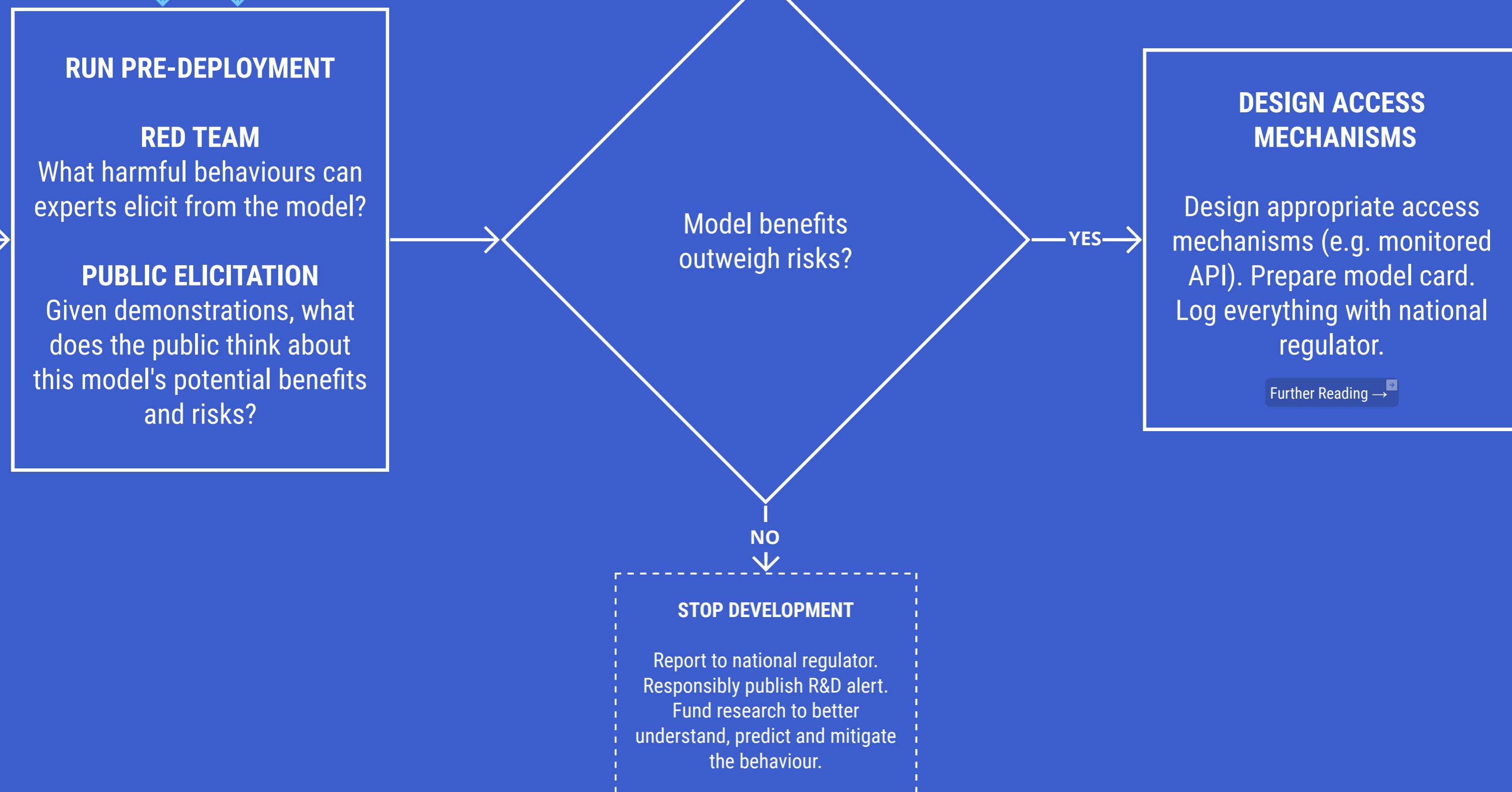
SPECIALISATION

Training for specific behaviours, goals or tasks



IV

PRE-DEPLOYMENT



EXCLUSIVITY PERIOD

PREPARE AND RELEASE REPORT

Prepare and release research paper / technical report while mindful of misuse potential

Further Reading →

ROLL OUT GRADUALLY

Monitor for misuse and unexpected behaviour.

SCALE UP ROLL-OUT

Monitor for misuse and unexpected behaviour.

ADDRESS ISSUES AS THEY COME UP

Share best practices with technical community and national regulator.

Following an exclusivity period of some years

MANAGE INCIDENTS

Establish incident reporting and redress mechanism, Establish adversarial testing and bounty programme.

Further Reading →

PREPARE LIABILITY REGIME

Prepare liability regime for harms caused by the model due to developer negligence.



