

Exploring Controllable Text Generation Techniques

Shrimai Prabhumoye, Alan W Black,
Ruslan Salakhutdinov



Carnegie Mellon University

Language Technologies Institute

Controllable Text Generation allows us to add “knobs” to control the attributes of the text to be generated

Controllable Text Generation allows us to add “knobs” to control the attributes of the text to be generated



**Who did
you like the best in
Avengers**

Controllable Text Generation allows us to add “knobs” to control the attributes of the text to be generated



Who did you like the best in Avengers



Control Attributes:
Robert Downey Jr.
Scarlett Johansson

Controllable Text Generation allows us to add “knobs” to control the attributes of the text to be generated



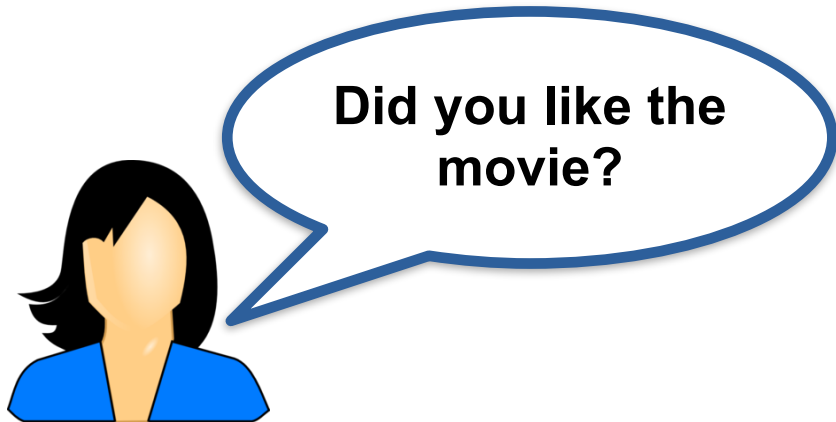
Who did you like the best in Avengers

I liked Robert Downey Jr. and Scarlett Johansson in the movie.

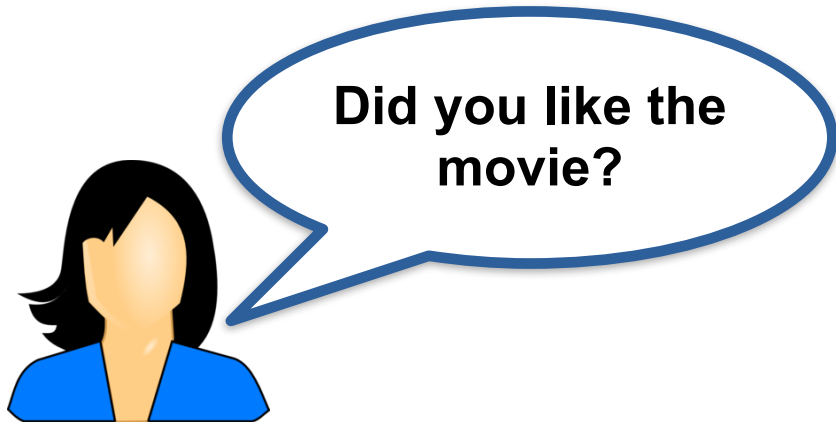


Control Attributes:
Robert Downey Jr.
Scarlett Johansson

Controllable Text Generation allows us to add “knobs” to control the attributes of the text to be generated



Controllable Text Generation allows us to add “knobs” to control the attributes of the text to be generated



**Control Attributes:
positive**

Controllable Text Generation allows us to add “knobs” to control the attributes of the text to be generated



Did you like the movie?



Yeah, I loved the movie!

**Control Attributes:
positive**

Controllable Text Generation allows us to add “knobs” to control the attributes of the text to be generated



Did you like the movie?



Yeah, I loved the movie!

Control Attributes:
negative

Controllable Text Generation allows us to add “knobs” to control the attributes of the text to be generated



Did you like the movie?

No, I hated it!



Control Attributes:
negative

Applications

- ***Dialogue System***
 - Persona, style of responses (polite, authority), content of responses, topic of conversation
- Recommend ***polite emails***
- ***Story Generation***
 - plot, ending, sentiment, topic, persona
- ***Report Generation*** (websites, Wikipedia articles)

Motivation

Style Transfer Through Back-Translation

Shrimai Prabhumoye, Yulia Tsvetkov, Ruslan Salakhutdinov, Alan W Black
Carnegie Mellon University, Pittsburgh, PA, USA
{sprabhum, ytsvetko, rsalakhu, awb}@cs.cmu.edu

Abstract

Style transfer is the task of rephrasing the text to contain specific stylistic properties without changing the intent or affect within the context. This paper introduces

These goals have motivated a considerable amount of recent research efforts focused at “controlled” language generation—aiming at separating the semantic content of *what* is said from the stylistic dimensions of *how* it is said. These include approaches relying on heuristic substitu-

Published as a conference paper at ICLR 2019

OF WIKIPEDIA: KNOWLEDGE-POWERED CONVERSATIONAL AGENTS

Emily Dinan*, Stephen Roller*, Kurt Shuster*, Angela Fan, Michael Auli, Jason Weston
Facebook AI Research
{edinan, roller, kshuster, angelafan, michaelauli, jase}@fb.com

ABSTRACT

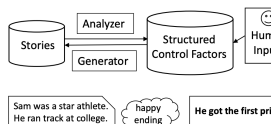
In open-domain dialogue intelligent agents should exhibit the use of knowledge, however there are few convincing demonstrations of this to date. The most popular sequence to sequence models typically “generate and hope” generic utterances

Towards Controllable Story Generation

Nanyun Peng, Marjan Ghazvininejad, Jonathan May, Kevin Knight
Information Sciences Institute & Computer Science Department
University of Southern California
{npeng, ghazvini, jonmay, knight}@isi.edu

Abstract

We present a general framework of analyzing existing story corpora to generate controllable and creative new stories. The proposed framework needs little manual annotation to achieve controllable story generation. It creates a new



Towards Content Transfer through Grounded Text Generation

Shrimai Prabhumoye, Chris Quirk, Michel Galley
Carnegie Mellon University, Microsoft Research
5000 Forbes Avenue, One Microsoft Way
Pittsburgh, PA 15219, Redmond, WA 98052
sprabhum@andrew.cmu.edu {chrisq, mgalley}@microsoft.com

Abstract

Recent work in neural generation has attracted significant interest in controlling the *form* of text, such as style, persona, and politeness. However, there has been less work on control-

Monkey selfie copyright dis

From Wikipedia, the free encyclopedia

The monkey selfie copyright dispute is a series of disputes about the



Politeness Transfer: A Tag and Generate Approach

Aman Madaan*, Amrith Setlur*, Tanmay Parekh*, Barnabas Poczos, Graham Neubig, Yiming Yang, Ruslan Salakhutdinov, Alan W Black, Shrimai Prabhumoye
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA, USA
{amadaan, asetlur, tparekh}@cs.cmu.edu

Abstract

This paper introduces a new task of politeness transfer which involves converting non-polite sentences to polite sentences while preserving the meaning. We also provide a dataset of more than 1.39 million instances automatically labeled for politeness to assess baseline

Rao and Tetreault, 2018; Xu et al., 2012; Jhantani et al., 2017) has not focused on politeness as a style transfer task, and we argue that defining it is cumbersome. While native speakers of a language and cohabitants of a region have a good working understanding of the phenomenon of politeness

The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)

Plan-and-Write: Towards Better Automatic Storytelling

Lili Yao,^{1,3*} Nanyun Peng,^{2*} Ralph Weischedel,² Kevin Knight,² Dongyan Zhao,¹ Rui Yan^{1†}
lilyao@tencent.com, {npeng, weisched, knight}@isi.edu
{zhaodongyan, ruiyan}@pku.edu.cn

¹Institute of Computer Science and Technology, Peking University
²Information Sciences Institute, University of Southern California, ³Tencent AI Lab

Abstract

Automatic storytelling is challenging since it requires generating long, coherent natural language to describe a sensible sequence of events. Despite considerable efforts on automatic story generation in the past, prior work either is restricted in

Title (Given)	The Bike Accident
Storyline (Extracted)	Carrie → bike → sneak → nervous → leg
Story (Human Written)	Carrie had just learned how to ride a bike. She didn't have a bike of her own. Carrie would sneak rides on her

Motivation

Style Transfer Through Back-Translation

Shrimai Prabhumoye, Yulia Tsvetkov, Ruslan Salakhutdinov, Alan W Black
Carnegie Mellon University, Pittsburgh, PA, USA
{sprabhum,ytsvetko,rsalakhu,awb}@cs.cmu.edu

Abstract

Style transfer is the task of rephrasing the text to contain specific stylistic properties without changing the intent or affect within the context. This paper introduces

These goals have motivated a considerable amount of recent research efforts focusing on “controlled” language generation—aiming at controlling the semantic content of *what* is said and the stylistic dimensions of *how* it is said. This paper includes approaches relying on heuristic

Large body of work

Towards Controllable Story Generation

Politeness Transfer: A Tag and Generate Approach

Aman Madaan*, Amrith Setlur*, Tanmay Parekh*, Barnabas Poczos, Graham Neubig, Shuangping Yang, Ruslan Salakhutdinov, Alan W Black, Shrimai Prabhumoye
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA, USA
{amadaan, asetlur, tparekh}@cs.cmu.edu

Abstract

This paper introduces a new task of politeness transfer, which involves converting non-polite sentences to polite sentences while preserving the original meaning. We also provide a dataset of 1.39 million instances automatically generated for politeness to aggression benchmark.

Rao and Tetreault, 2018; Xu et al., 2012; Jhamtani et al., 2017) has not focused on politeness as a style transfer task, and we argue that defining it is cumbersome. While native speakers of a language and cohabitants of a region have a good working understanding of the phenomenon of politeness

Published as a conference paper at ICLR 2019

Towards Content Transfer through Grounded Text Generation

The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)

OF WIKIPEDIA: KNOWLEDGE-POWERED CONVERSATIONAL AGENTS

Emily Dinan*, Stephen Roller*, Kurt Shuster*, Angela Fan, Michael Auli, Jason Weston
Facebook AI Research
{edinan,roller,kshuster,angelifan,michaelauli,jase}@fb.com

Shrimai Prabhumoye
Carnegie Mellon University
5000 Forbes Avenue
Pittsburgh, PA 15219
sprabhum@andrew.cmu.edu

Chris Quirk, Michel Galley
Microsoft Research
One Microsoft Way
Redmond, WA 98052
{chrisq,mgalley}@microsoft.com

Plan-and-Write: Towards Better Automatic Storytelling

Lili Yao,^{1,3*} Nanyun Peng,^{2*} Ralph Weischedel,² Kevin Knight,² Dongyan Zhao,¹ Rui Yan^{1†}
lilyao@tencent.com, {npeng,weisched,knight}@isi.edu
{zhaodongyan,ruiyan}@pku.edu.cn

¹Institute of Computer Science and Technology, Peking University

²Information Sciences Institute, University of Southern California, ³Tencent AI Lab

Abstract

Recent work in neural generation has attracted significant interest in controlling the *form* of text, such as style, persona, and politeness. However, there has been less work on controlling

Monkey selfie copyright dis

From Wikipedia, the free encyclopedia

The monkey selfie copyright dispute is a series of disputes about the



Abstract

Automatic storytelling is challenging since it requires generating long, coherent natural language to describe a sensible sequence of events. Despite considerable efforts on automatic story generation in the past, prior work either is restricted in

ABSTRACT

In open-domain dialogue intelligent agents should exhibit the use of knowledge, however there are few convincing demonstrations of this to date. The most popular sequence to sequence models typically “generate and hope” generic utterances

Title (Given)	The Bike Accident
Storyline (Extracted)	Carrie → bike → sneak → nervous → leg
Story (Human Written)	Carrie had just learned how to ride a bike. She didn't have a bike of her own. Carrie would sneak rides on her

Motivation

Style Transfer Through Back-Translation

Shrimai Prabhumoye, Yulia Tsvetkov, Ruslan Salakhutdinov, Alan W Black
Carnegie Mellon University, Pittsburgh, PA, USA
{sprabhum,ytsvetko,rsalakhu,awb}@cs.cmu.edu

Abstract

Style transfer is the task of rephrasing the text to contain specific stylistic properties without changing the intent or affect within the context. This paper introduces

These goals have motivated a considerable amount of recent research efforts focusing on “controlled” language generation—aiming at controlling the semantic content of *what* is said while preserving the stylistic dimensions of *how* it is said. This paper includes approaches relying on heuristic

Large body of work

Towards Controllable Story Generation

Politeness Transfer: A Tag and Generate Approach

Aman Madaan*, Amrith Setlur*, Tanmay Parekh*, Barnabas Poczos, Graham Neubig, Shuangping Yang, Ruslan Salakhutdinov, Alan W Black, Shrimai Prabhumoye
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA, USA
{amadaan, asetlur, tparekh}@cs.cmu.edu

Abstract

This paper introduces a new task of politeness transfer, which involves converting non-polite sentences into polite sentences while preserving the original meaning. We also provide a dataset of 1.39 million instances automatically generated for politeness to encourage benchmark

Rao and Tetreault, 2018; Xu et al., 2012; Jhamtani et al., 2017) has not focused on politeness as a style transfer task, and we argue that defining it is cumbersome. While native speakers of a language and cohabitants of a region have a good working understanding of the phenomenon of politeness

Published as a conference paper at ICLR 2019

Towards Content Transfer through Grounded Text Generation

The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)

WIKIPEDIA: KNOWLEDGE-POWERED CONVERSATIONAL AGENTS

Emily Dinan*, Stephen Roller*, Kurt Shuster*, Angela Fan, Michael Auli, Jason Weston
Facebook AI Research
{edinan,roller,kshuster,angelaFan,michaelauli,jasonw}@fb.com

ABSTRACT

In open-domain dialogue intelligent agents should exhibit the use of knowledge, however there are few convincing demonstrations of this to date. The most popular sequence to sequence models typically “generate and hope” generic utterances

No unifying theme

Plan-and-Write: Towards Programmatic Storytelling

Nanyun Peng,^{2*} Ralph Weichert,¹ Dongyan Zhao,¹ Rui Yan^{1†}
lilyao@tencent.com, rweichert@isi.edu, zhaodongyan@tencent.com, ryan@isi.edu
¹Institute of Computer Science, University of California, San Diego
²Information Sciences Institute, University of California, Los Angeles
[†]Tencent AI Lab

Abstract

Programmatic storytelling is challenging since it requires generating natural language to describe a sequence of events. Despite considerable effort in the past, prior work on programmatic storytelling has not learned how to ride a bike → sneak → nervous → ... in't have a bike of her ... could sneak rides on her

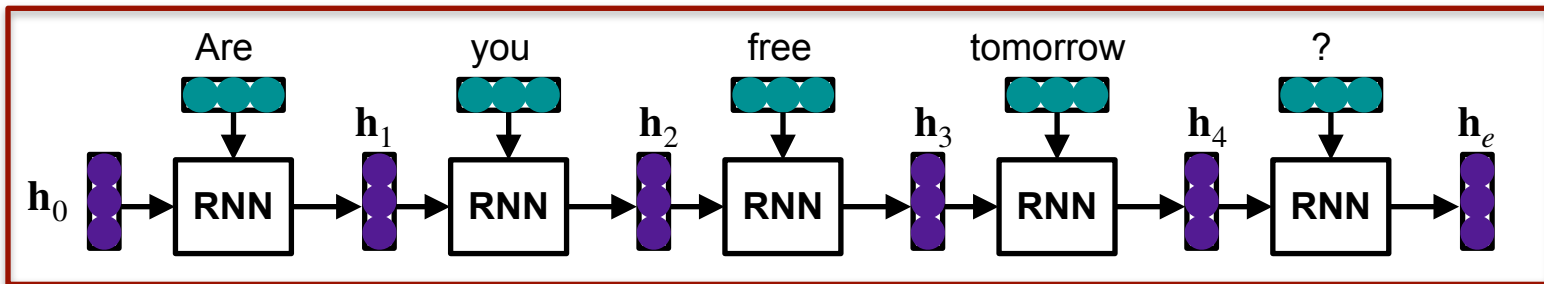


Contribution

- ***Controlled Generation Schema*** connects prior work
 - organize prior work
 - Schema contains 5 modules
 - Identify any architecture as belonging to one of these modules
 - Schema can be used with any algorithmic paradigm
- ***Collate knowledge*** about different techniques
 - Insights into the advantages of techniques
 - Pave way for new architectures
 - Provide easy access to comparison

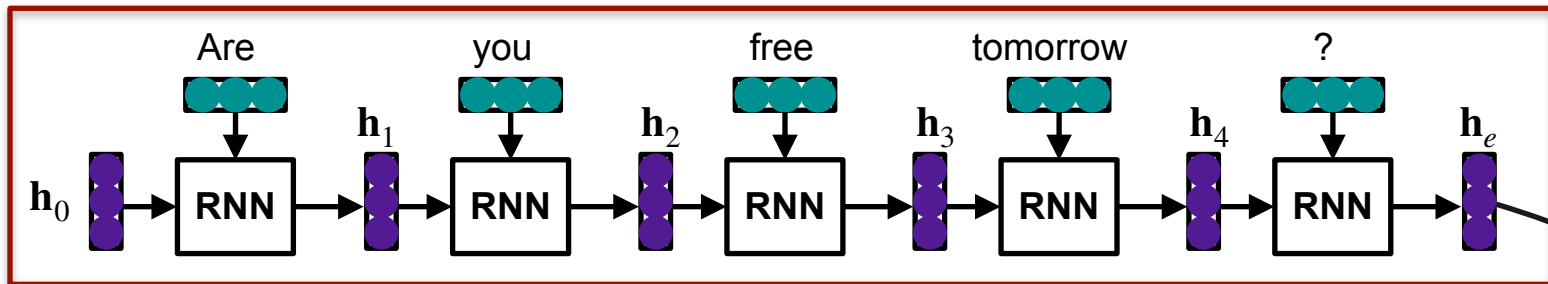
Generation Process

Encoder

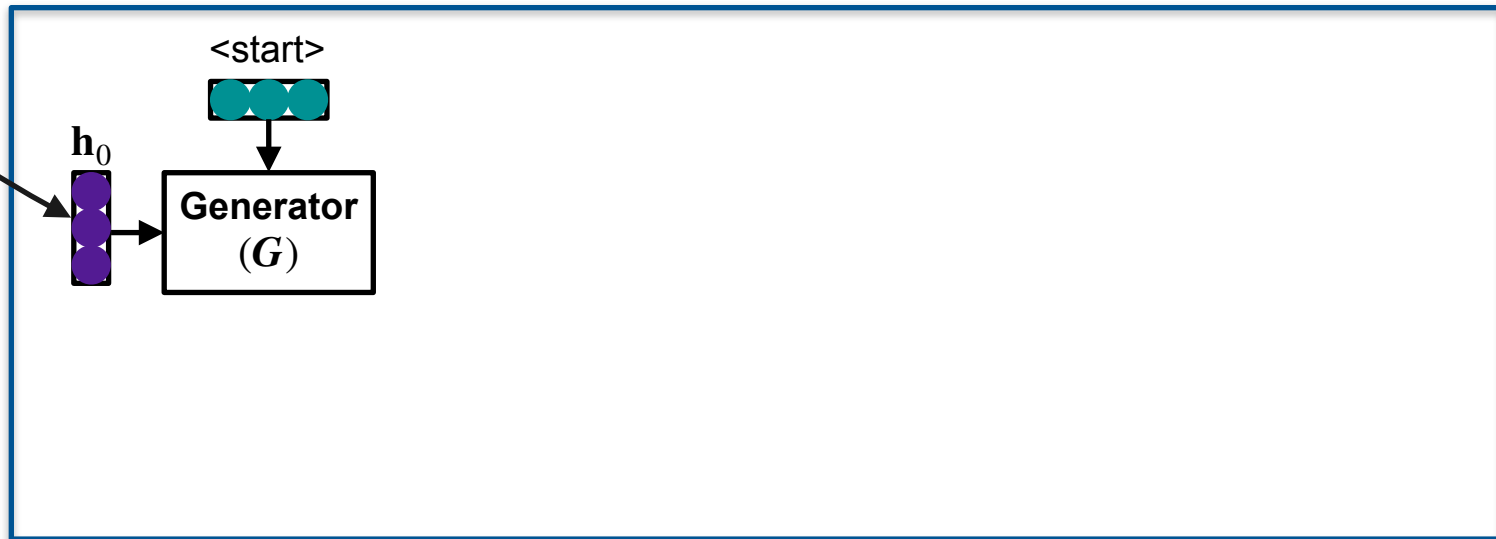


Generation Process

Encoder

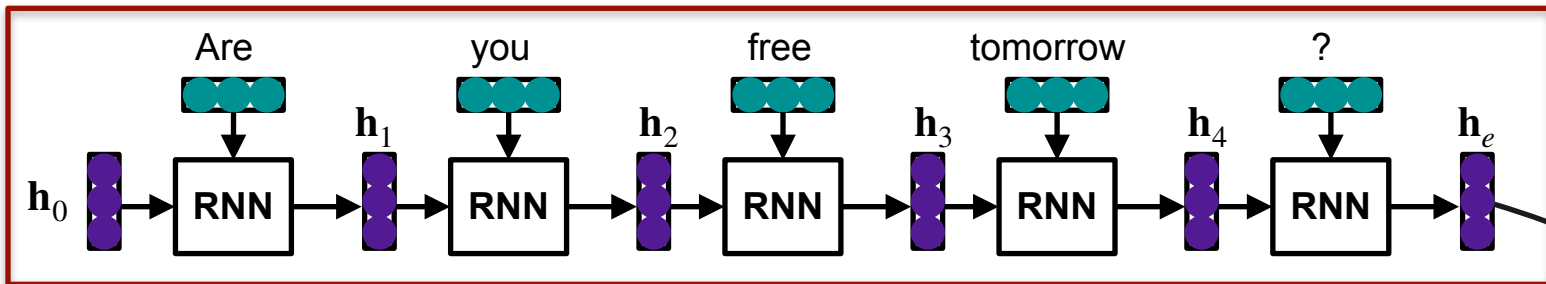


Decoder



Generation Process

Encoder

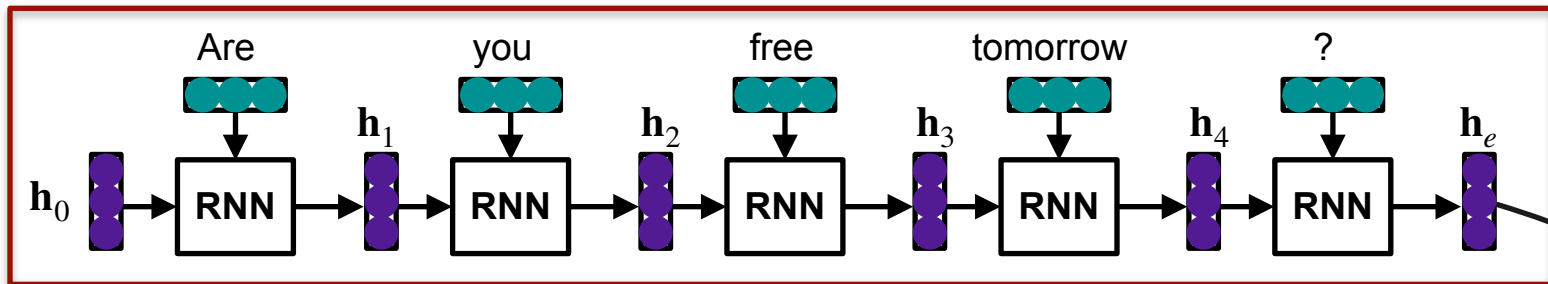


Decoder

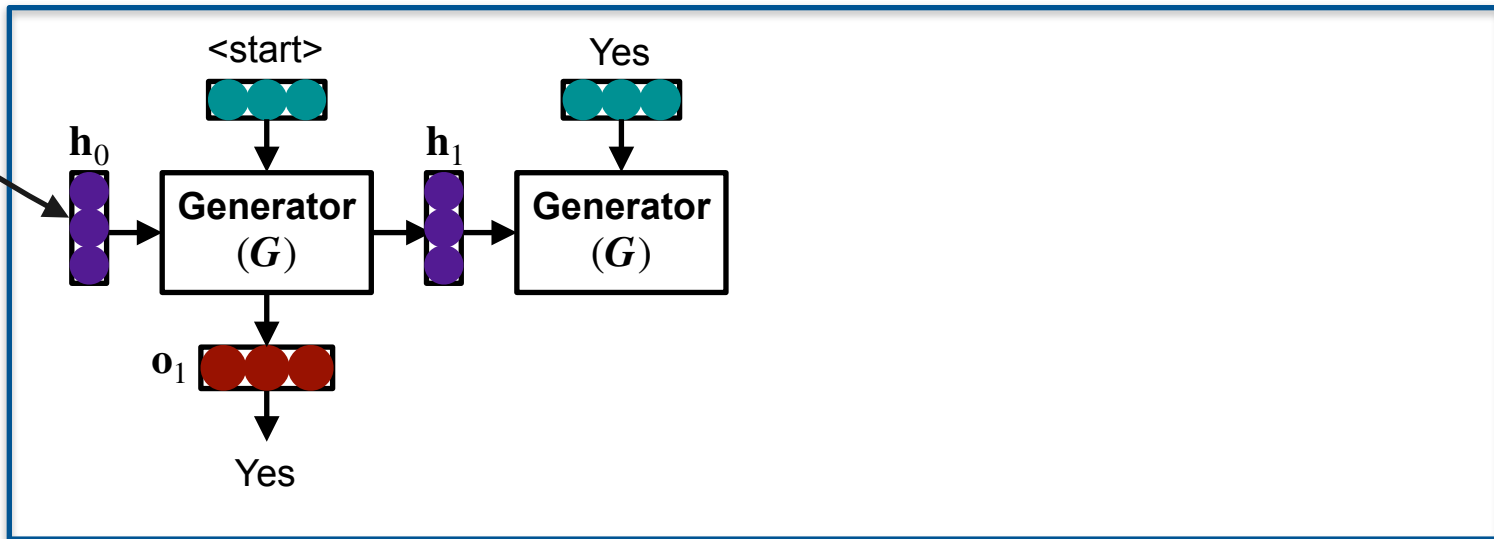


Generation Process

Encoder

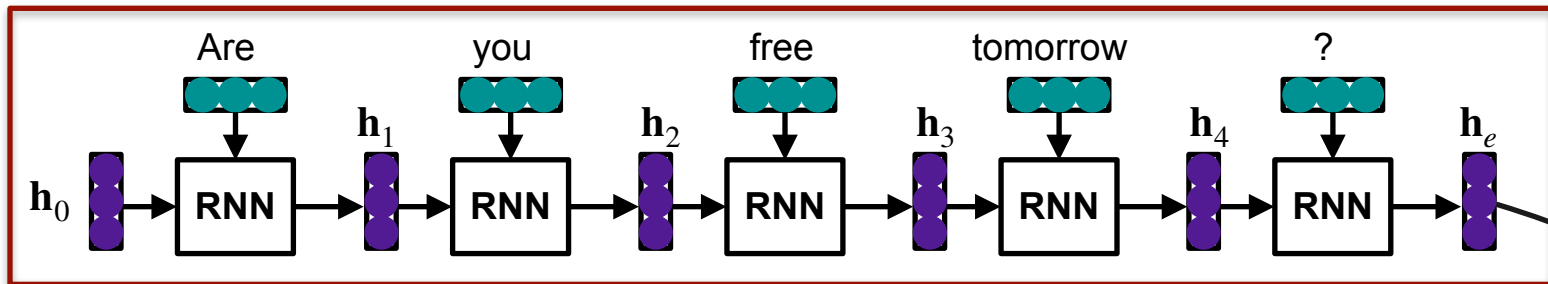


Decoder

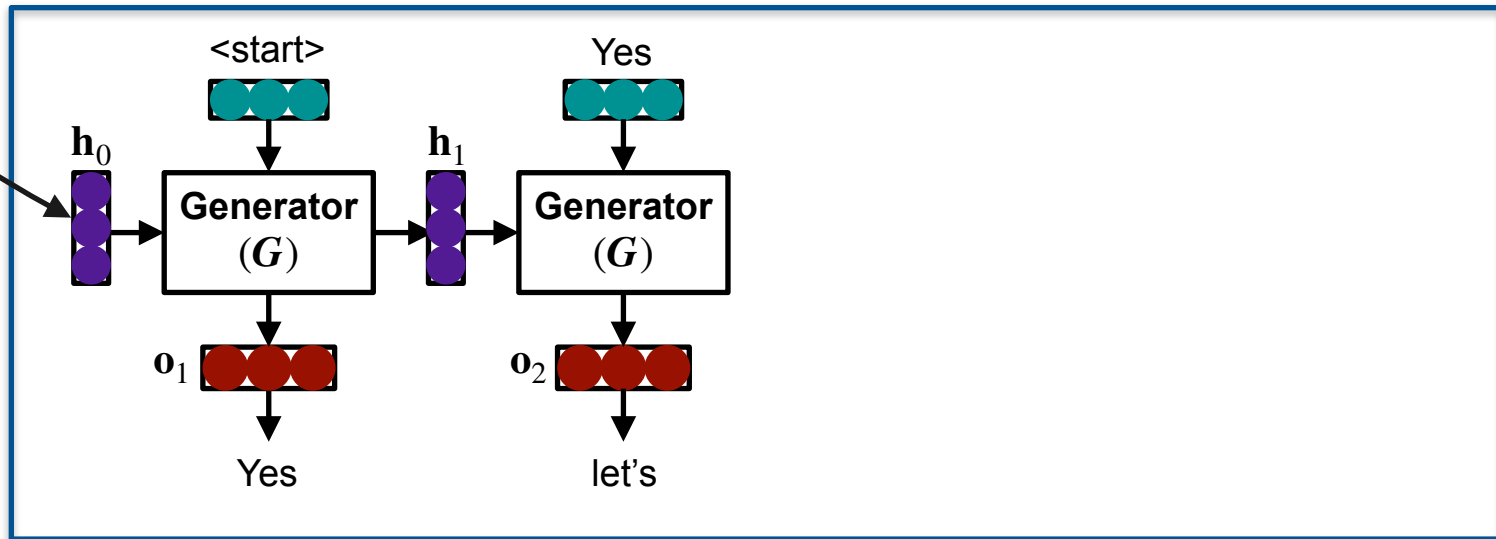


Generation Process

Encoder

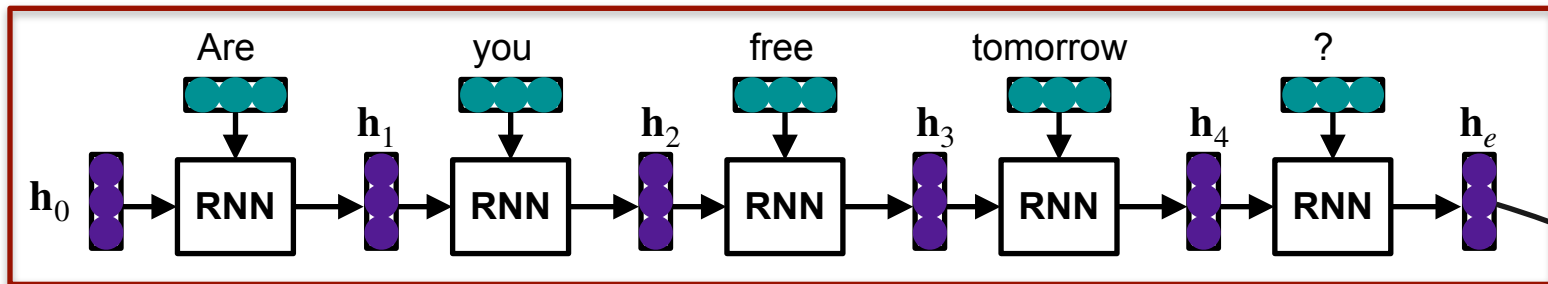


Decoder

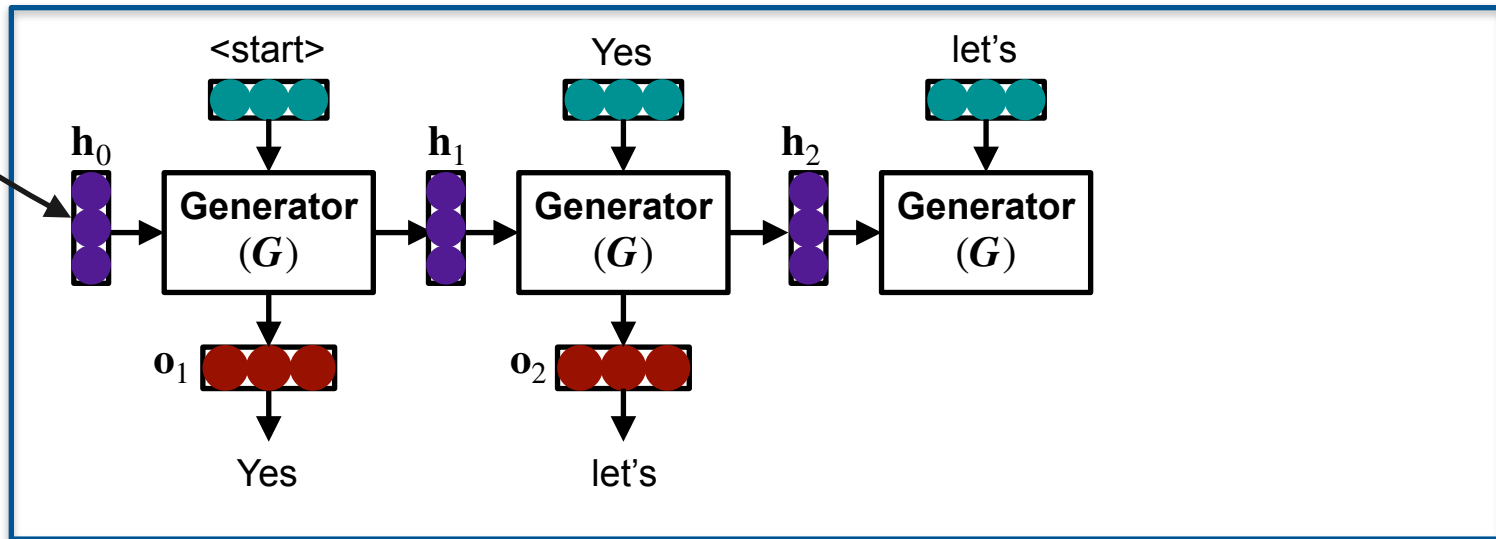


Generation Process

Encoder

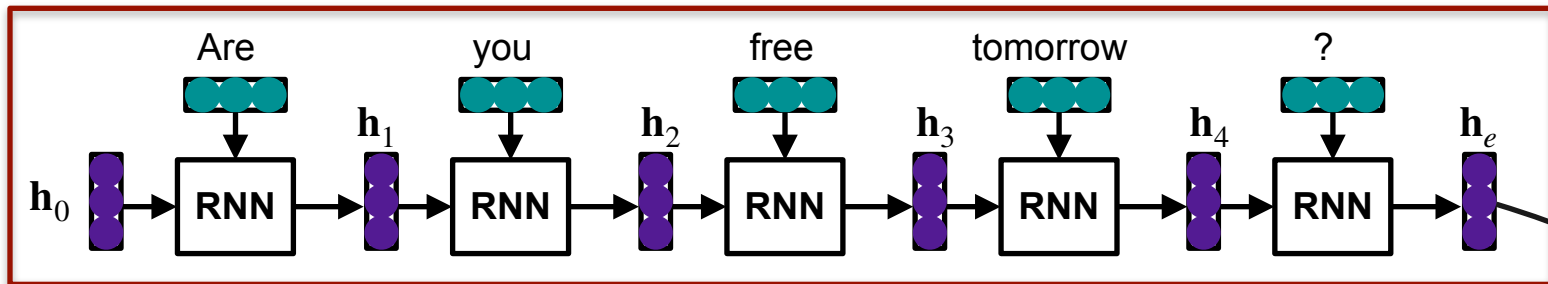


Decoder

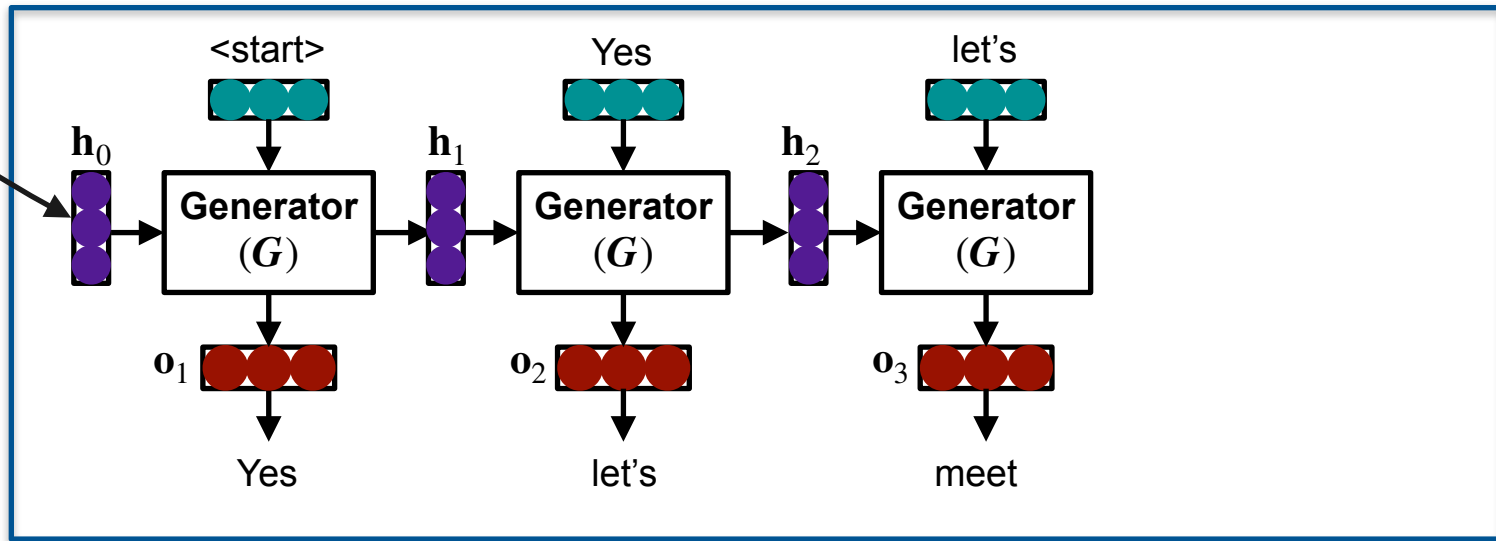


Generation Process

Encoder

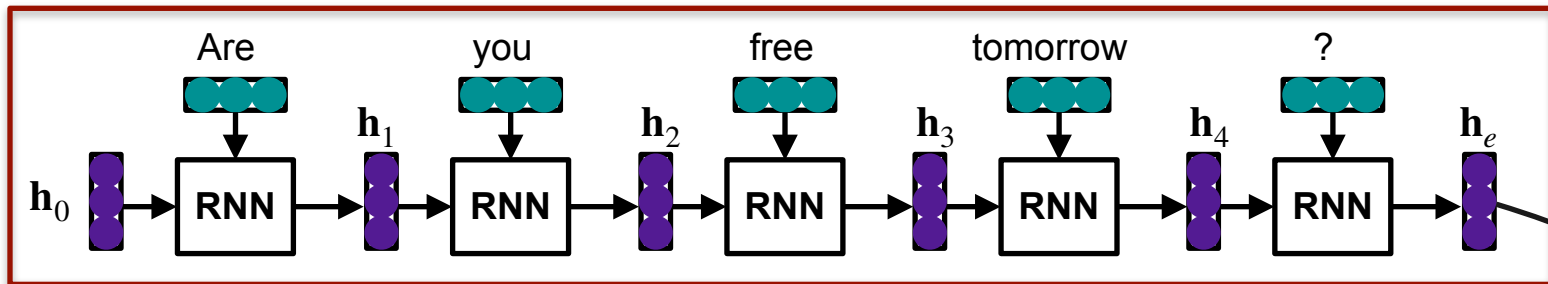


Decoder

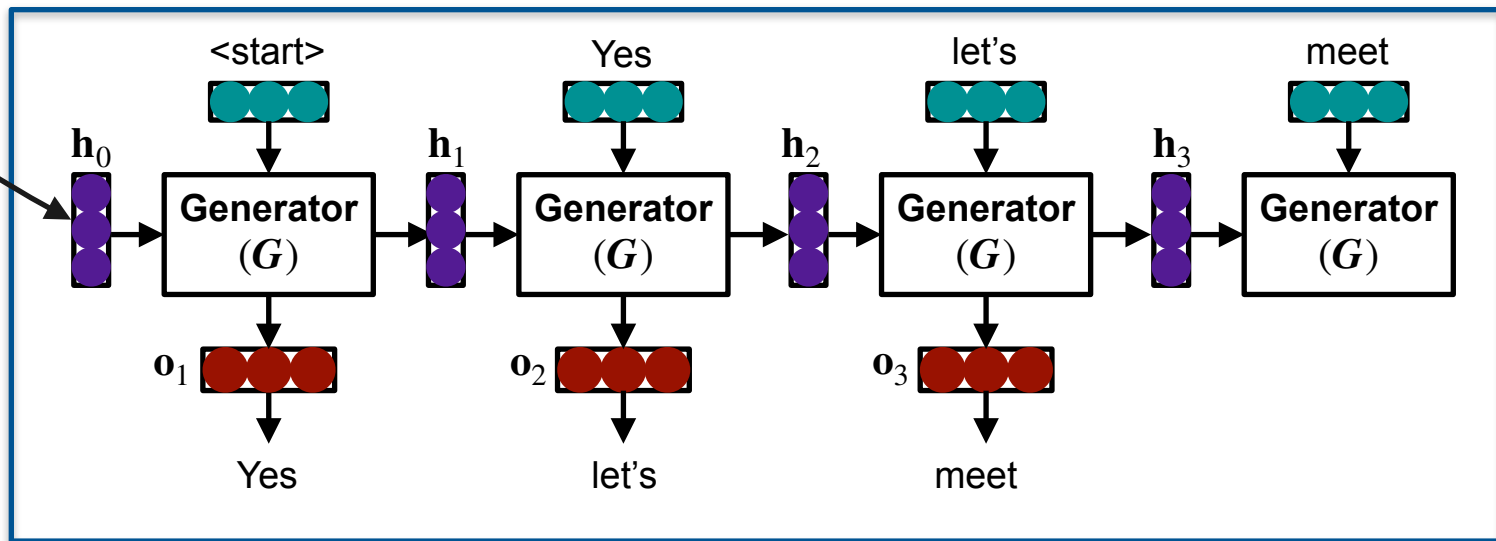


Generation Process

Encoder

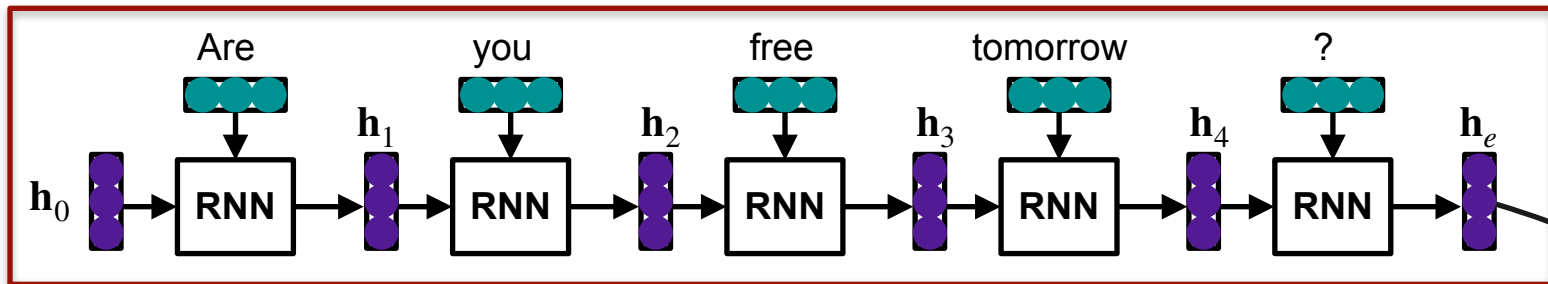


Decoder

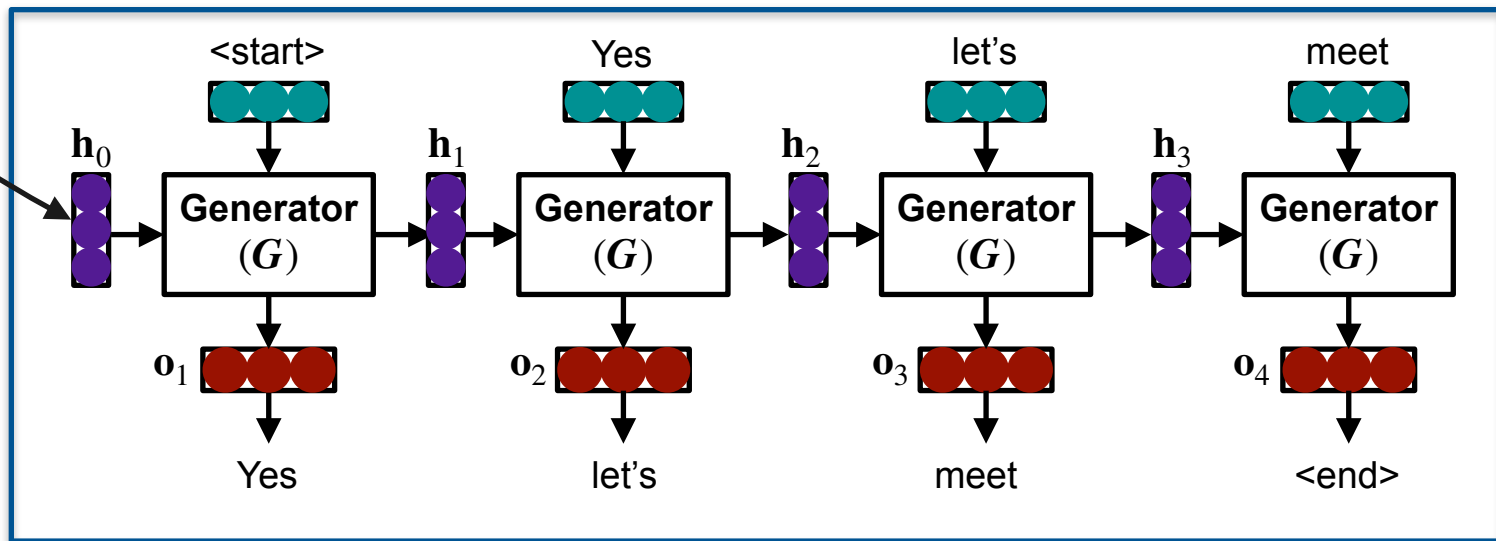


Generation Process

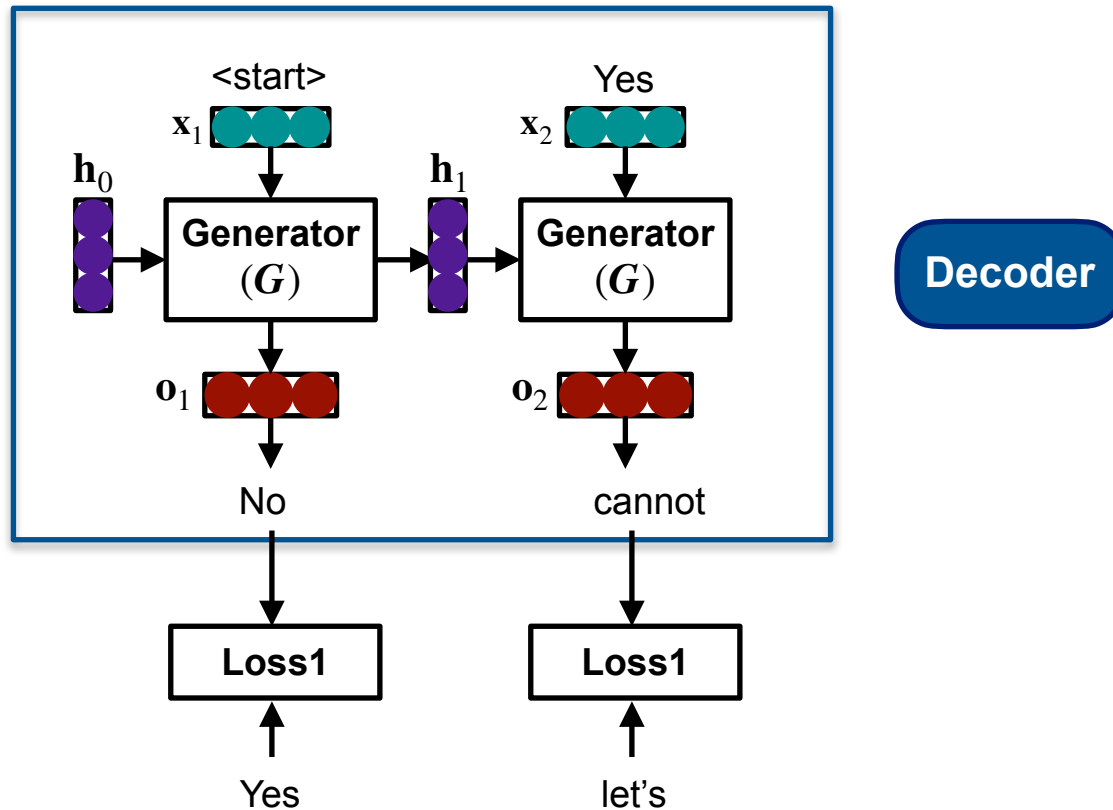
Encoder



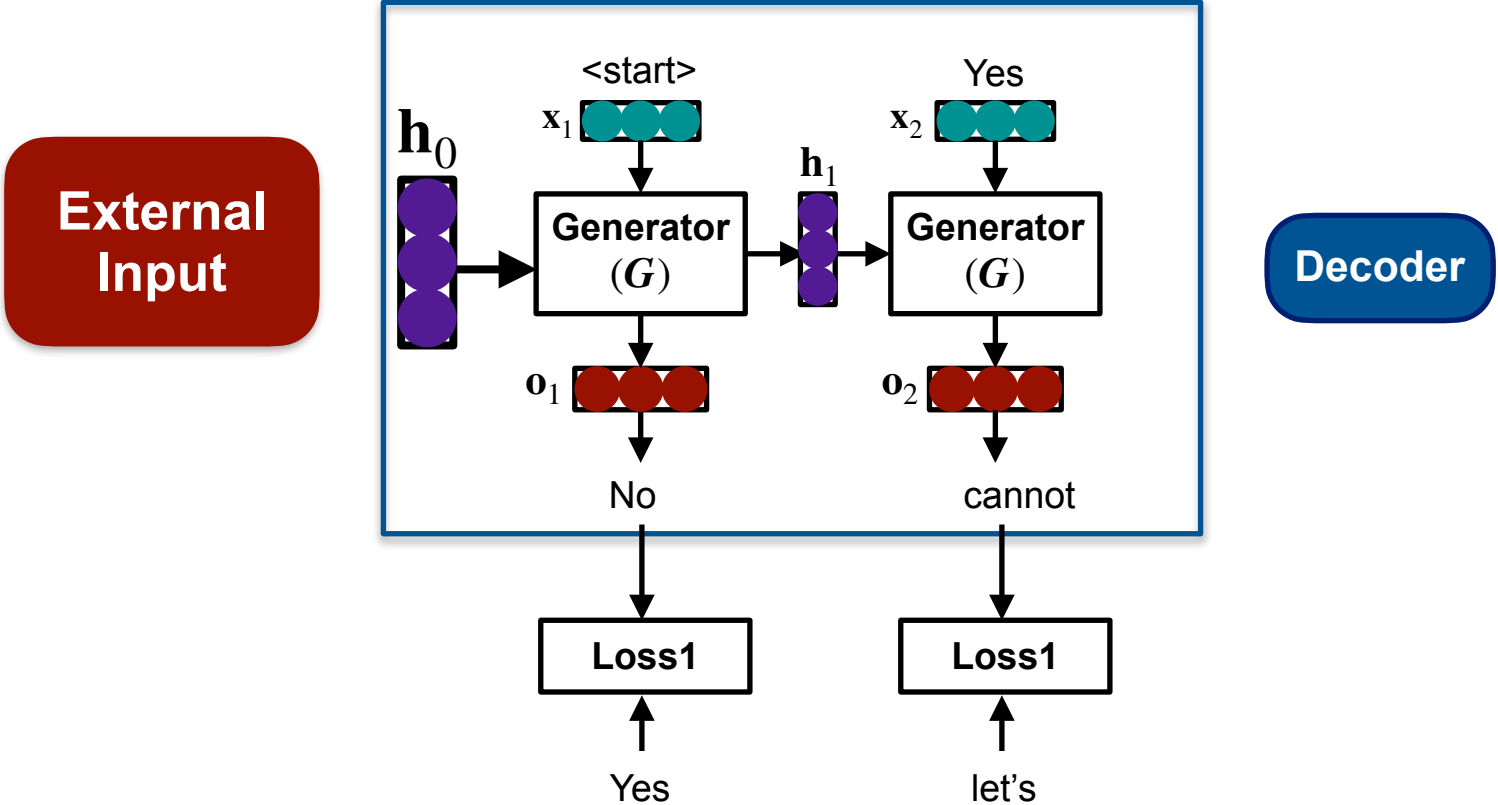
Decoder



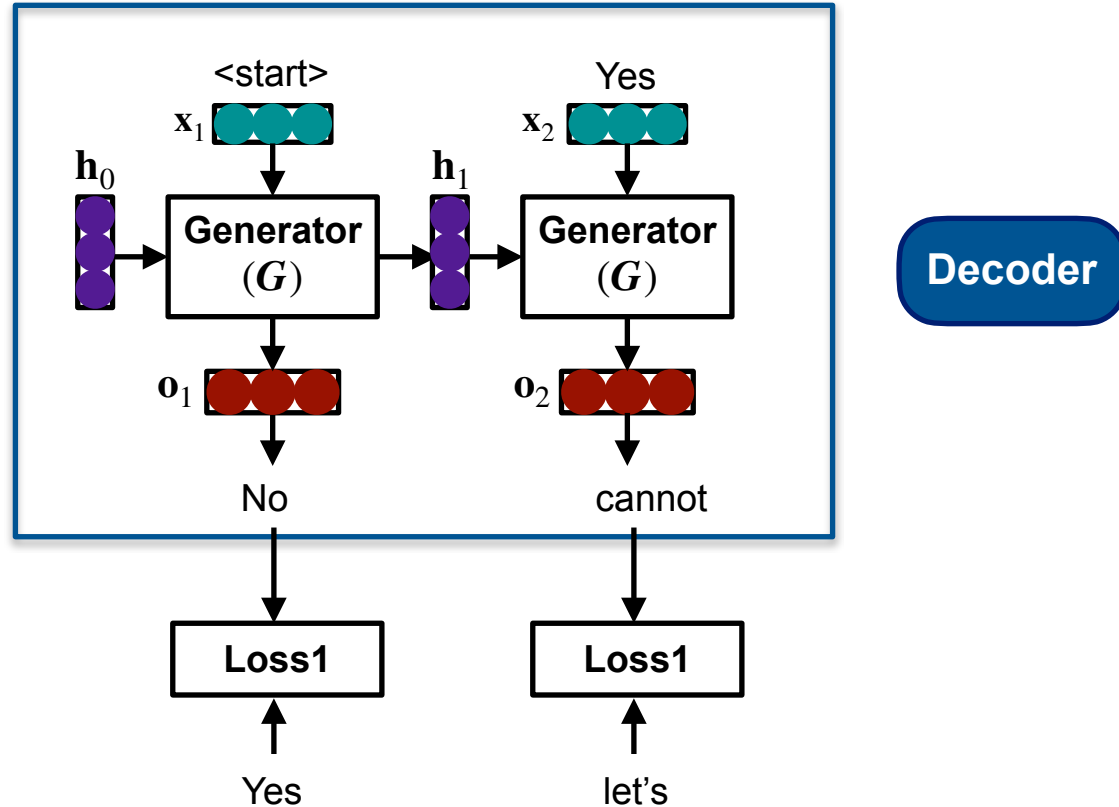
Modification Space



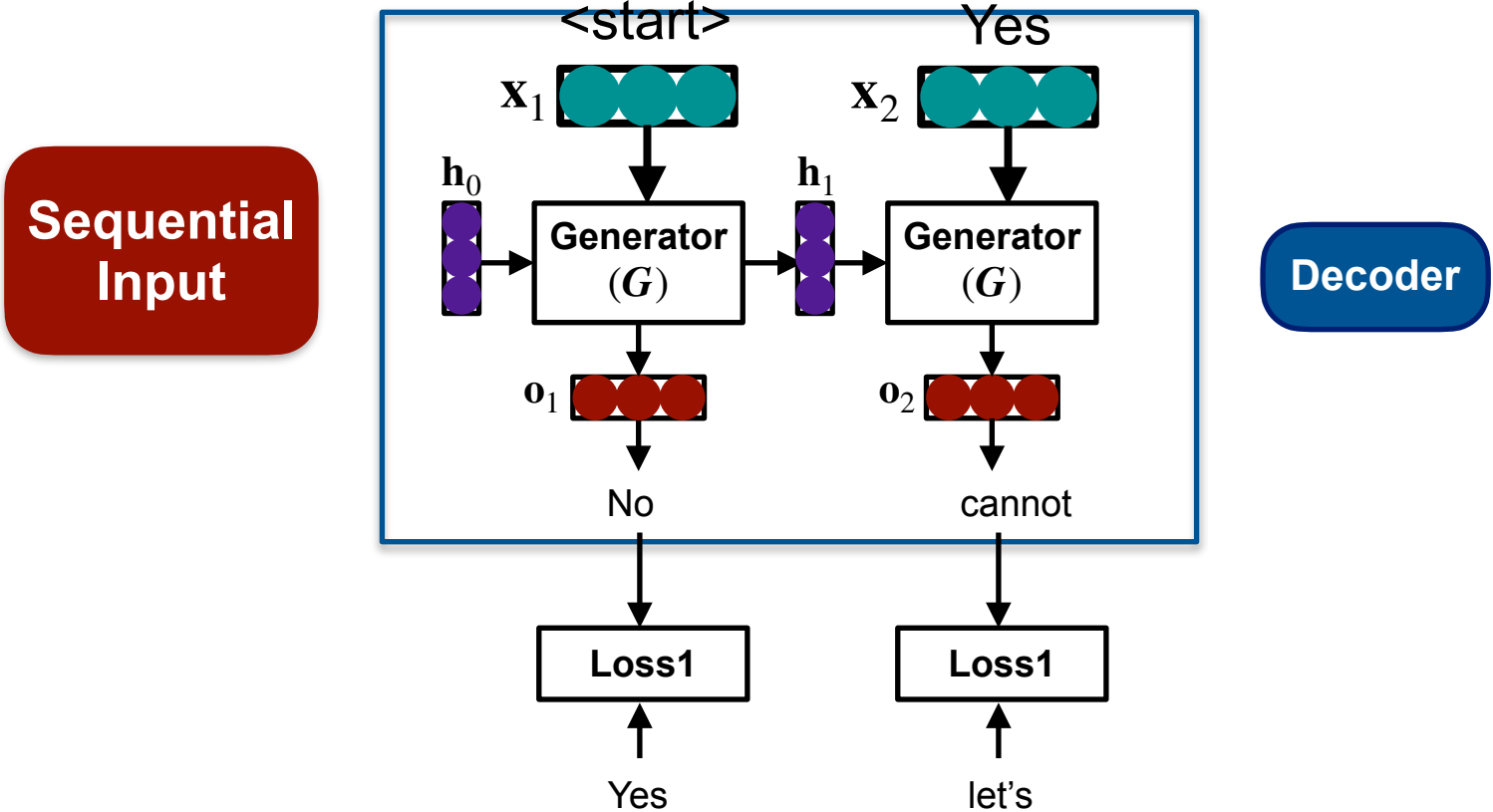
Modification Space



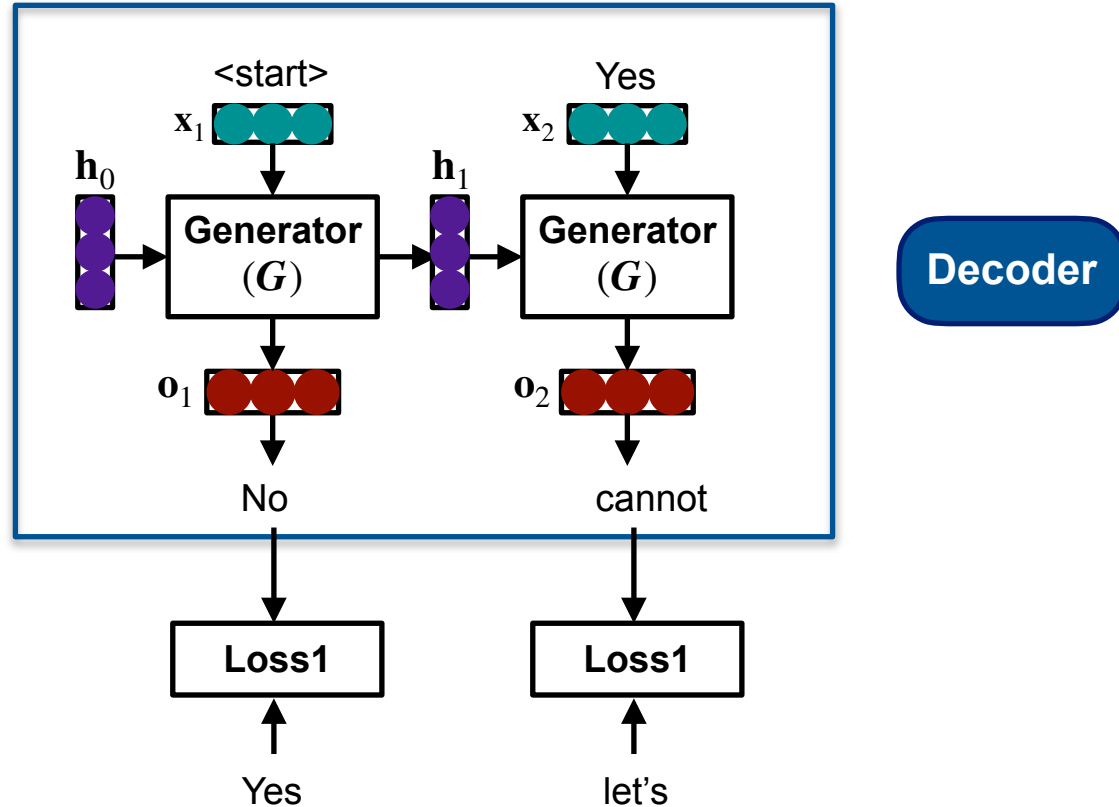
Modification Space



Modification Space

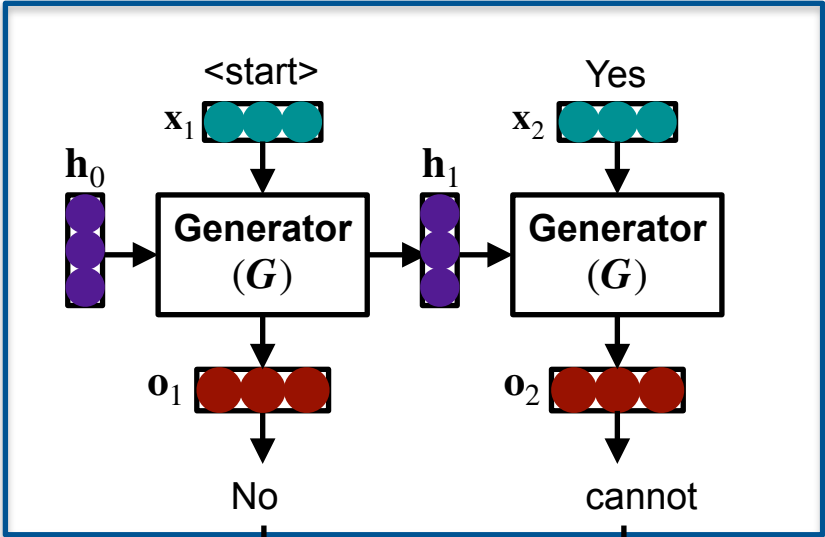


Modification Space



Modification Space

Generator



Decoder

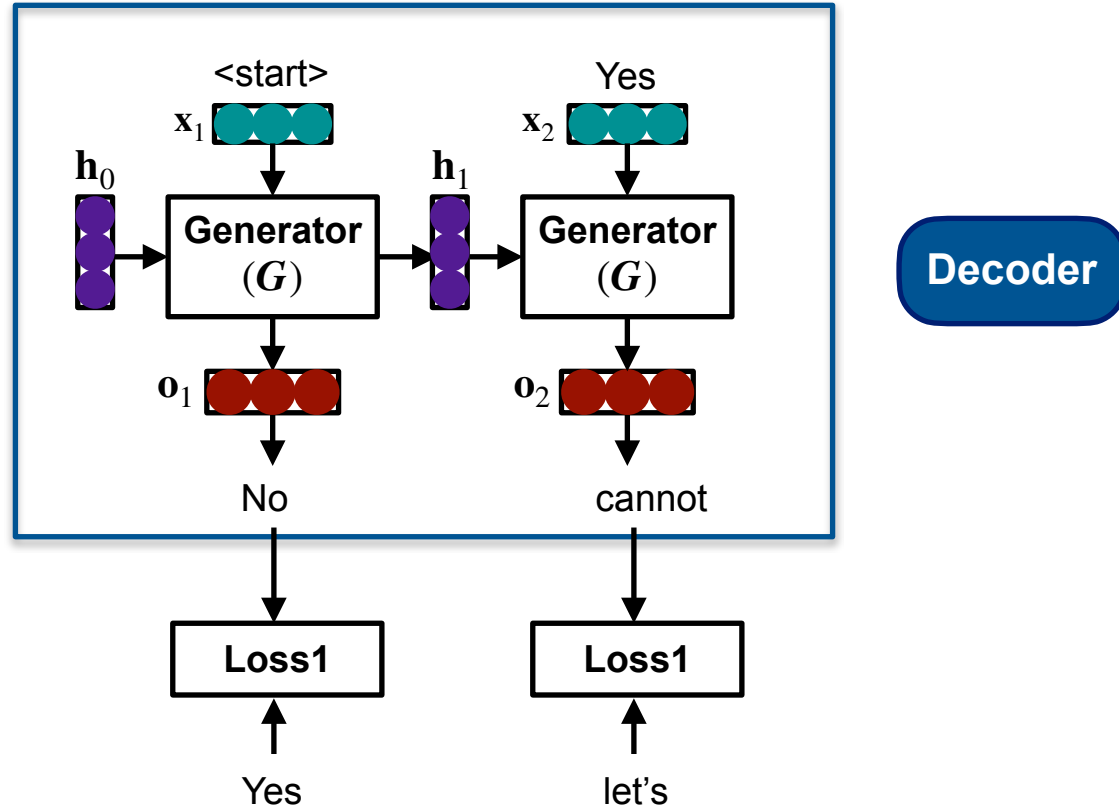
Loss1

Yes

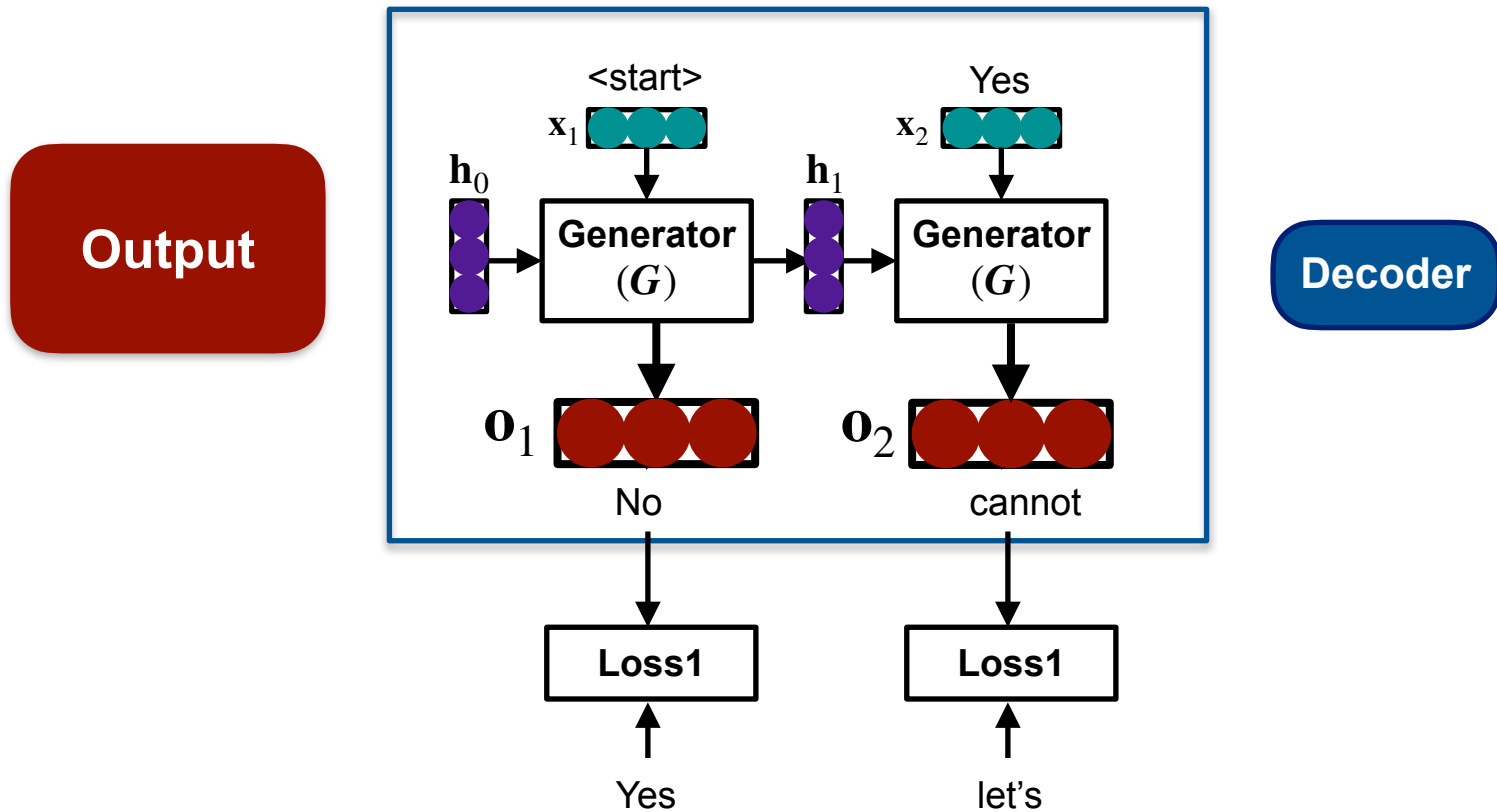
Loss1

let's

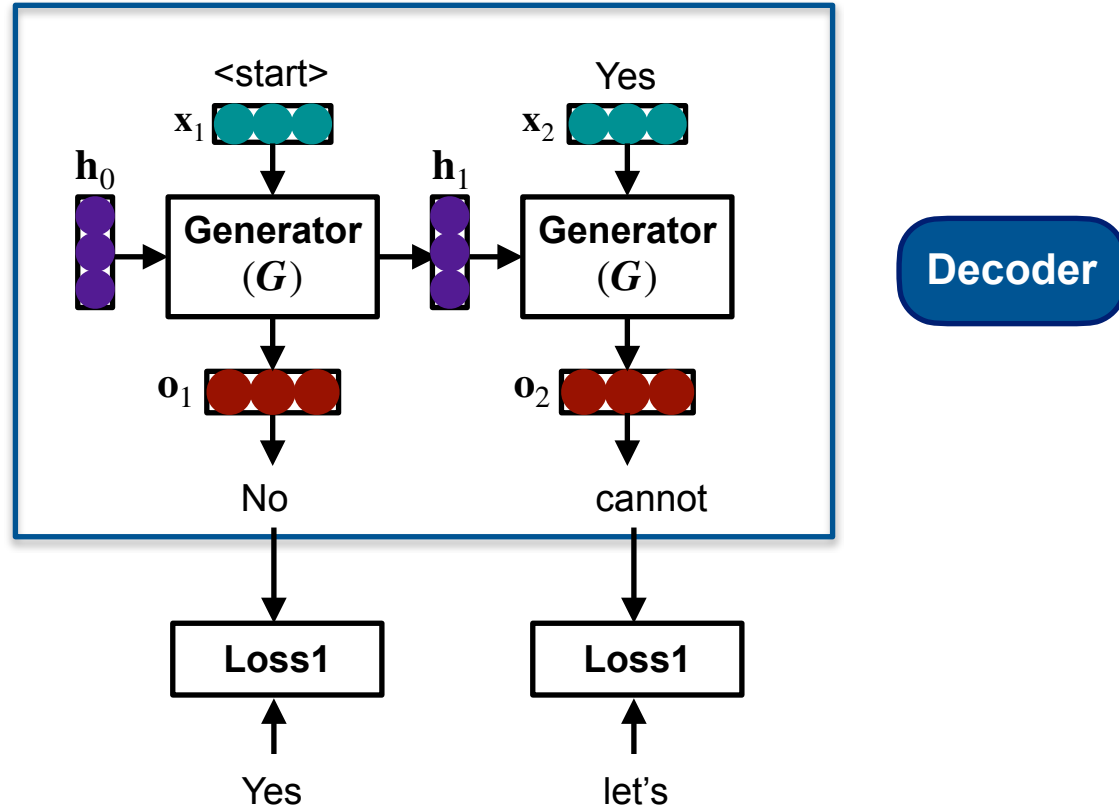
Modification Space



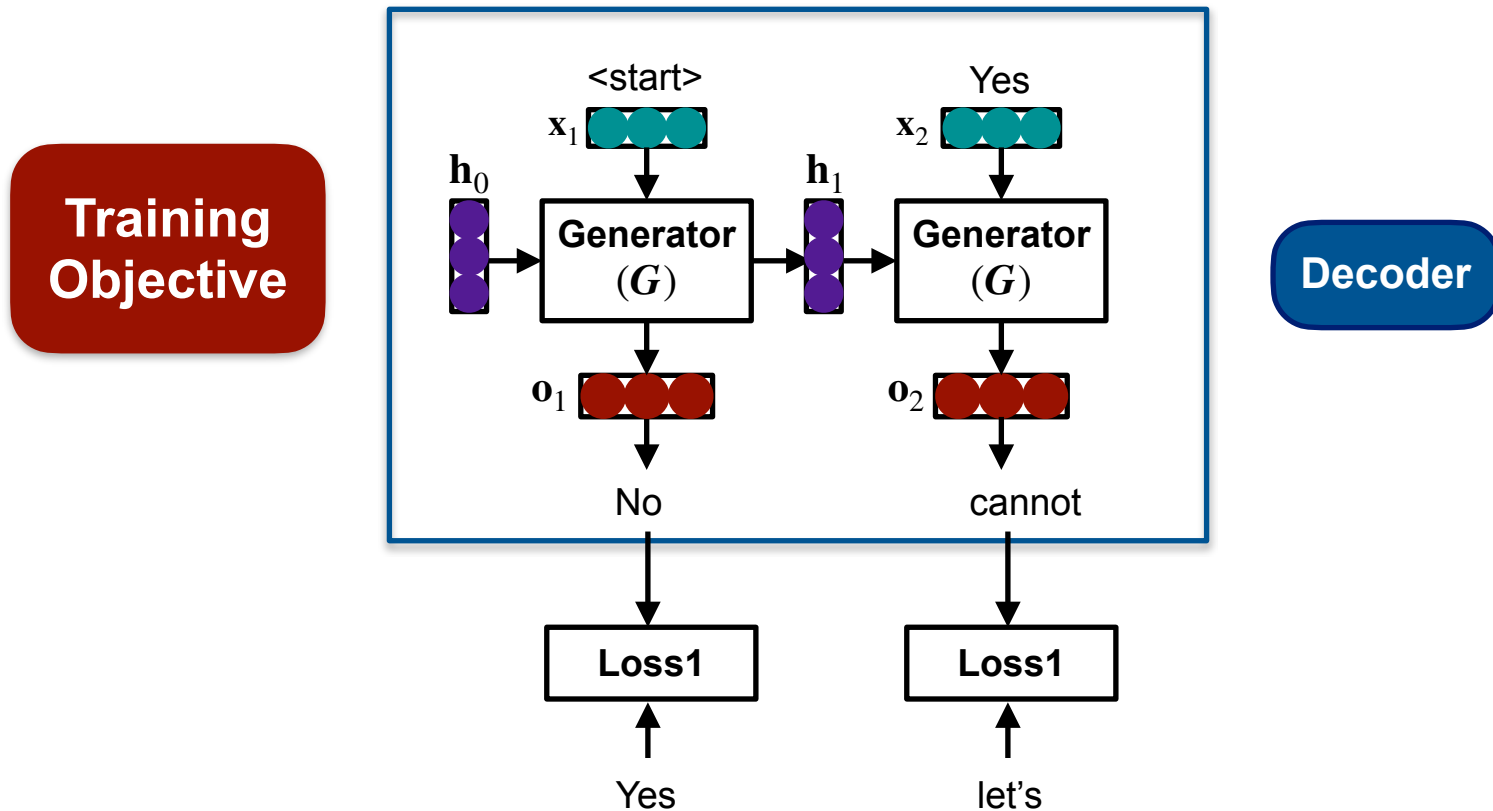
Modification Space



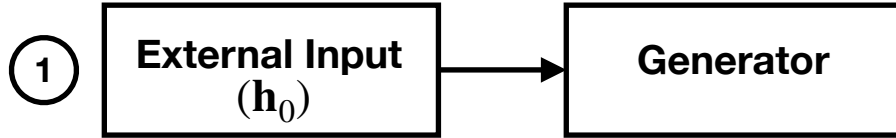
Modification Space



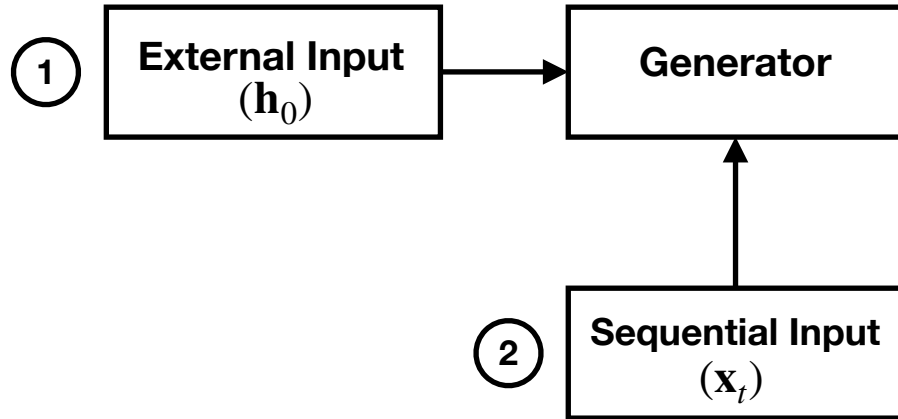
Modification Space



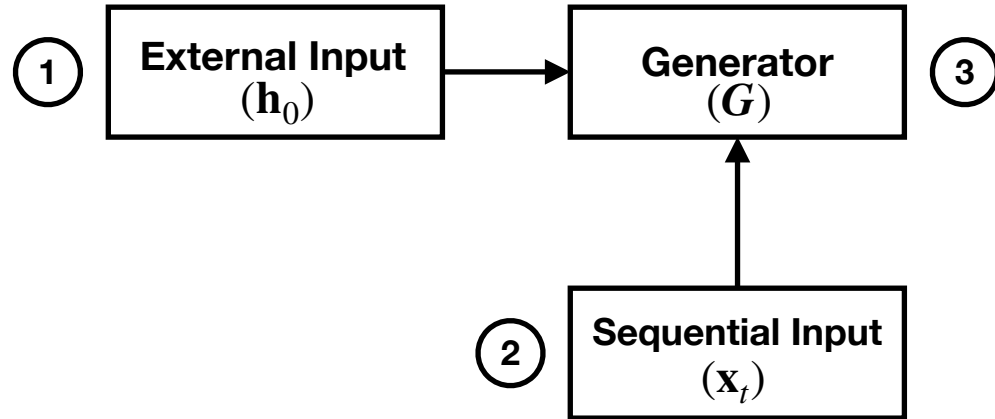
The proposed Schema



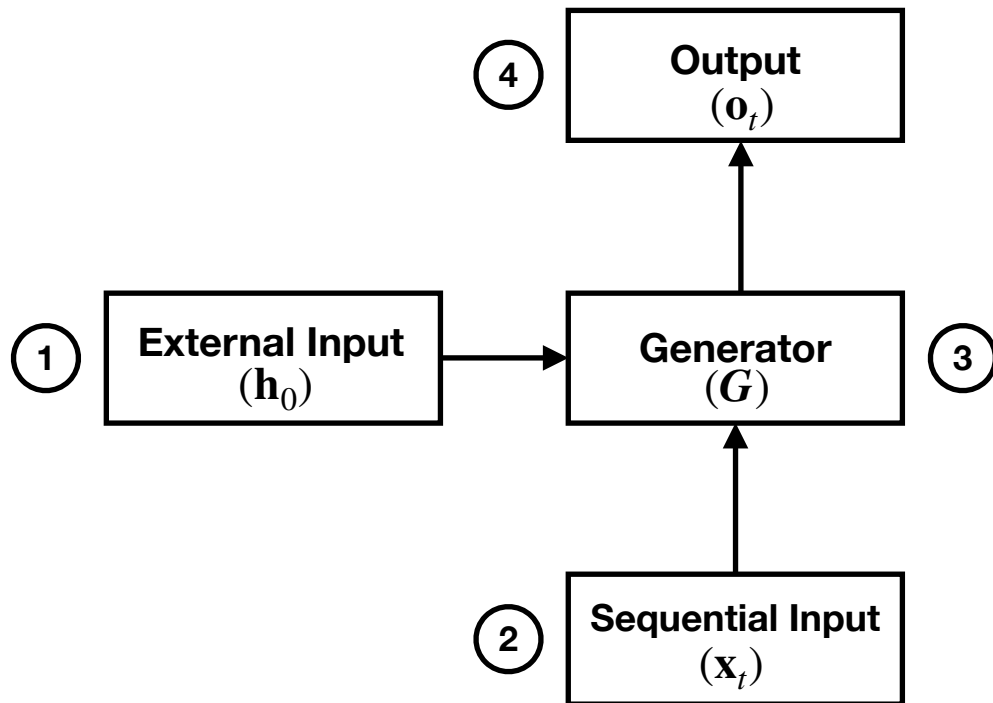
The proposed Schema



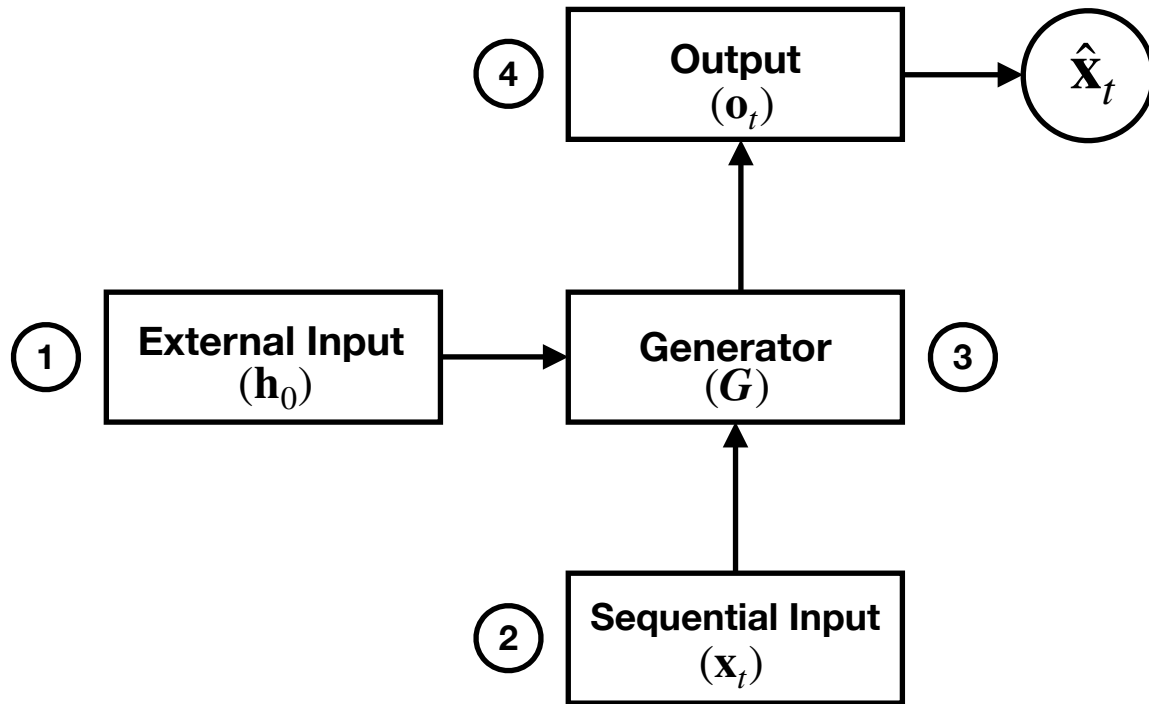
The proposed Schema



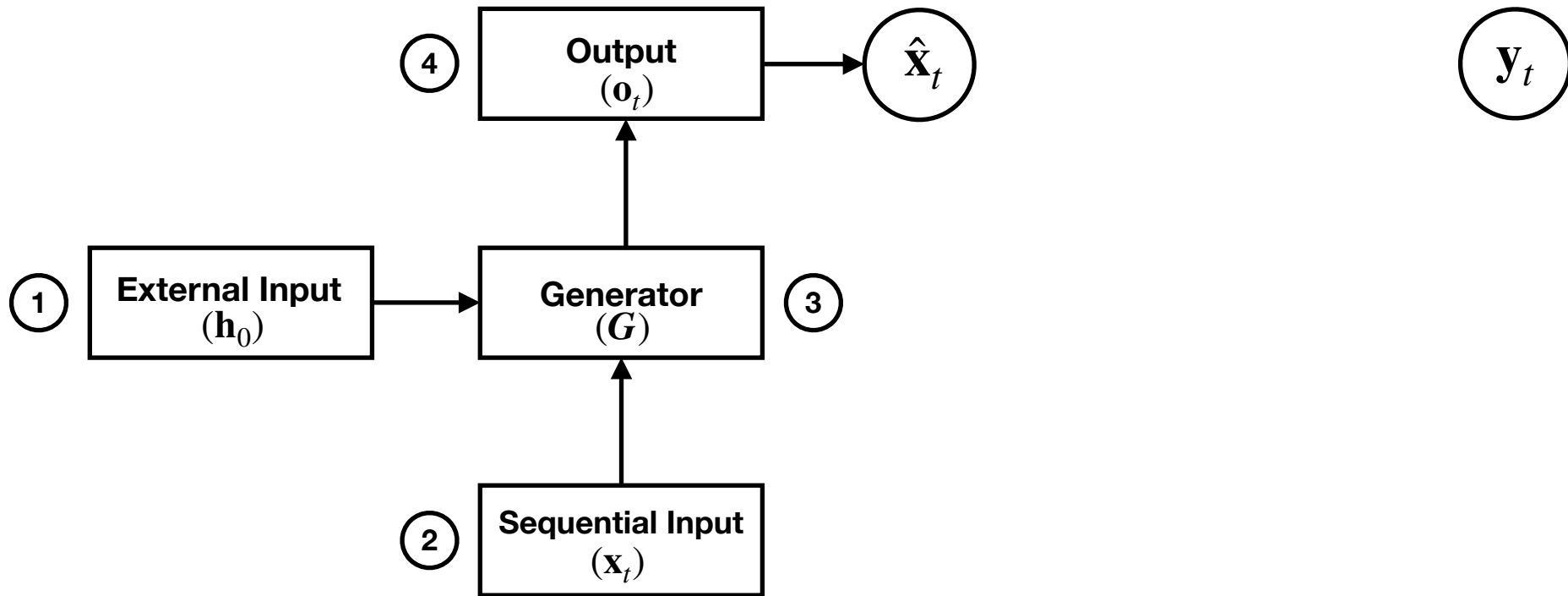
The proposed Schema



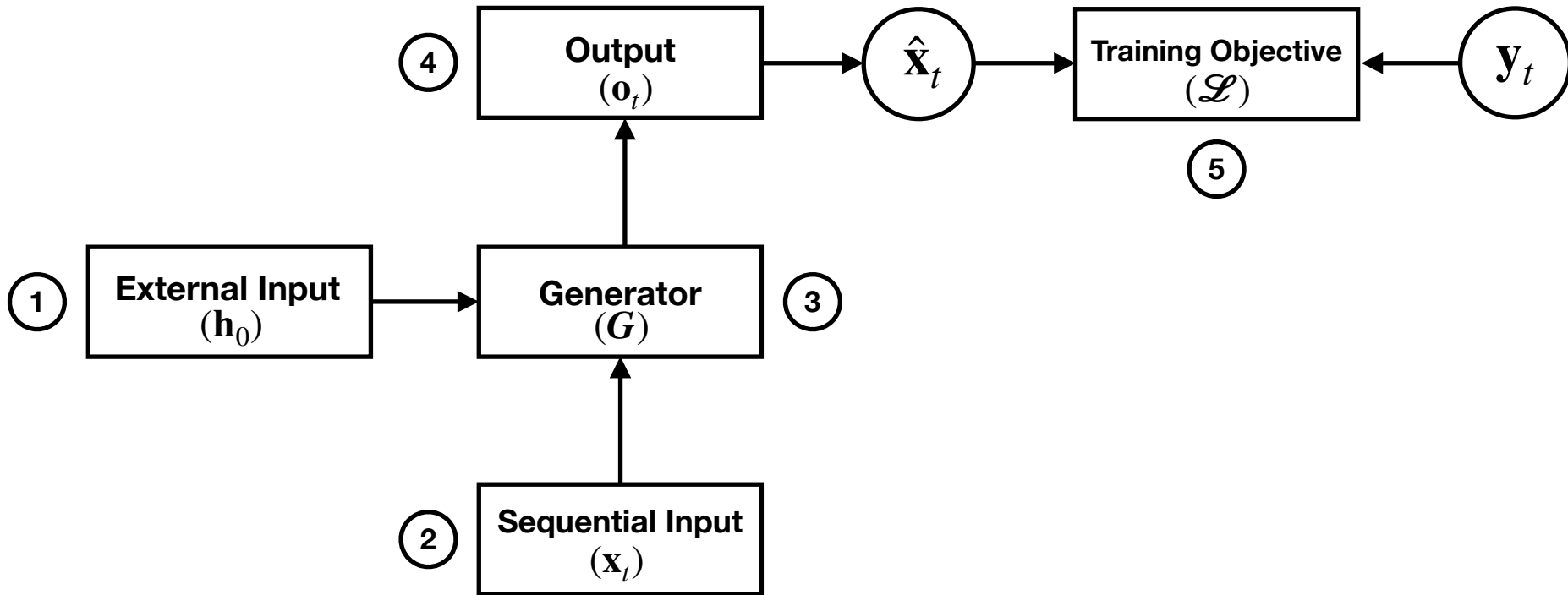
The proposed Schema



The proposed Schema



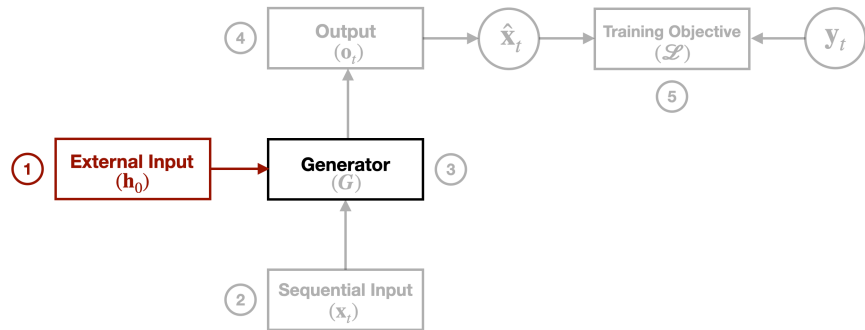
The proposed Schema



External Input

1. Decompose

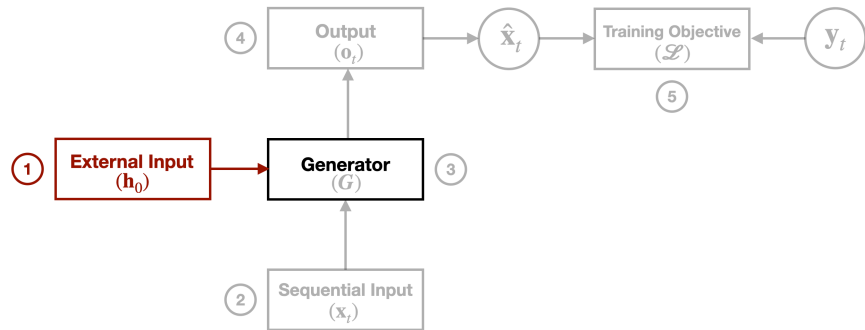
- \mathbf{h}_e decomposed into subspaces
- Provides *interpretable* representations
- Input should contain signal of control attribute
- Supervision on decomposed space



External Input

1. Decompose

- \mathbf{h}_e decomposed into subspaces
- Provides *interpretable* representations
- Input should contain signal of control attribute
- Supervision on decomposed space

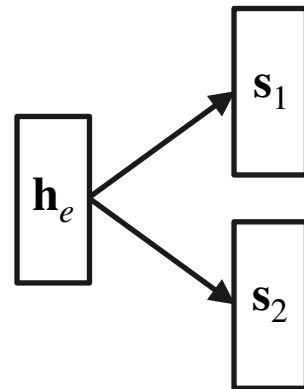
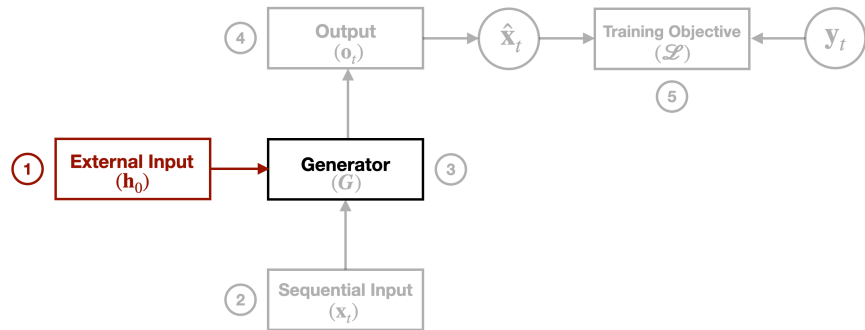


\mathbf{h}_e

External Input

1. Decompose

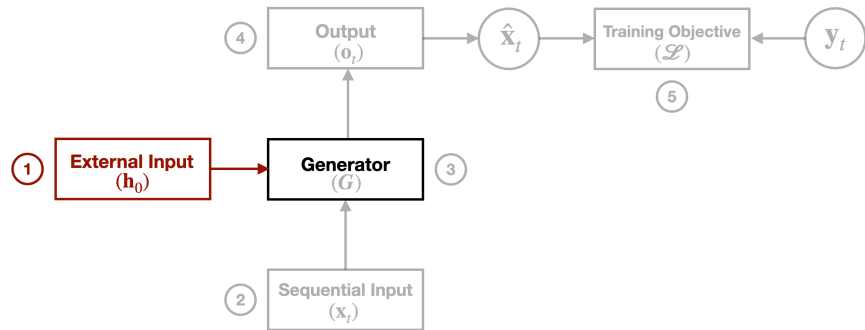
- \mathbf{h}_e decomposed into subspaces
- Provides *interpretable* representations
- Input should contain signal of control attribute
- Supervision on decomposed space



External Input

1. Decompose

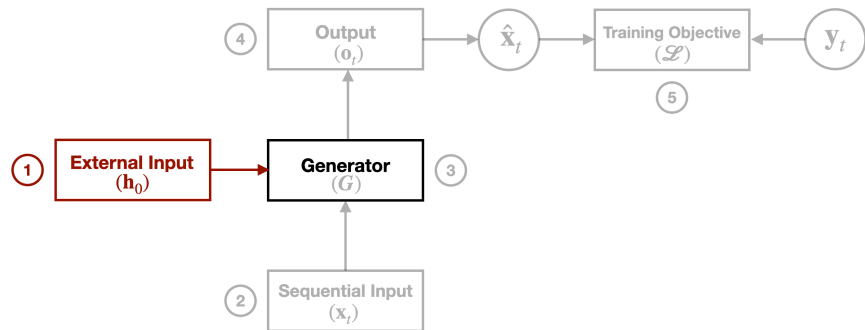
- \mathbf{h}_e decomposed into subspaces
- Provides *interpretable* representations
- Input should contain signal of control attribute
- Supervision on decomposed space



External Input

2. External Feedback

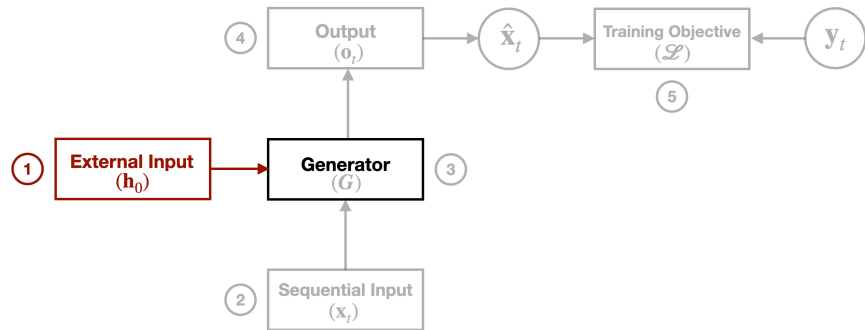
- regularizer to control \mathbf{h}_e
- must be jointly trained
- can be useful with decompose technique



External Input

2. External Feedback

- regularizer to control \mathbf{h}_e
- must be jointly trained
- can be useful with decompose technique

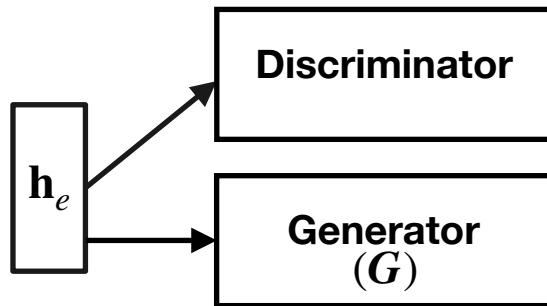
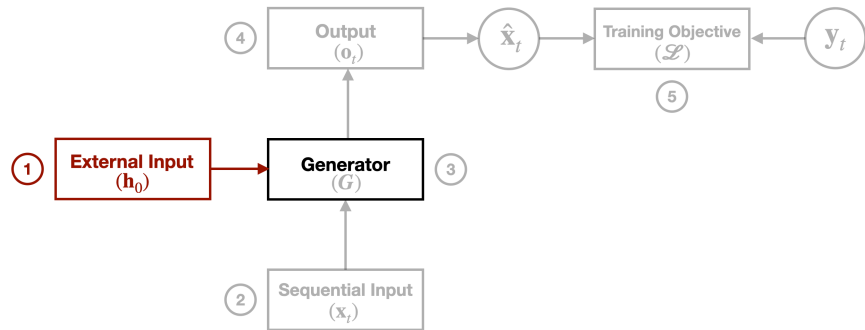


\mathbf{h}_e

External Input

2. External Feedback

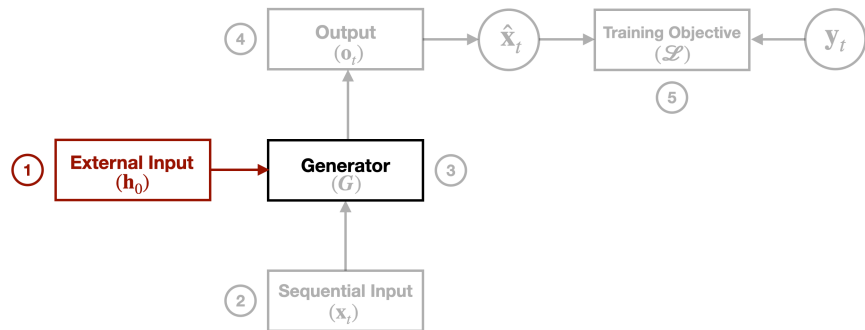
- regularizer to control \mathbf{h}_e
- must be jointly trained
- can be useful with decompose technique



External Input

2. External Feedback

- regularizer to control \mathbf{h}_e
- must be jointly trained
- can be useful with decompose technique



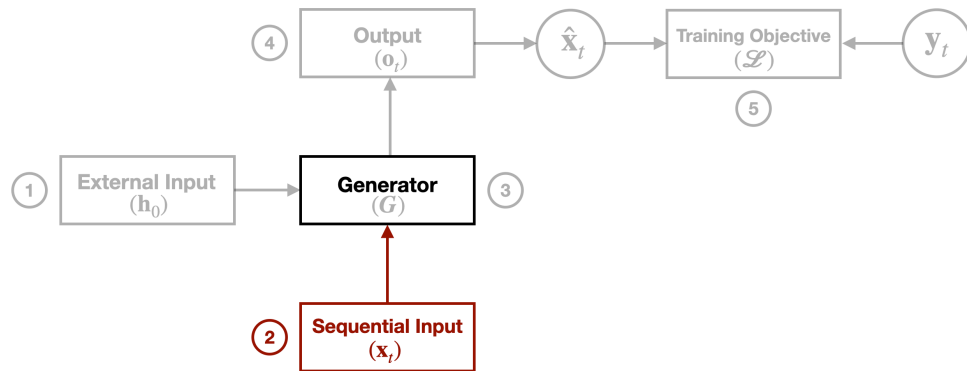
3. Arithmetic or Linear Transform

4. Stochastic Changes

Sequential Input

1. Arithmetic or Linear Transform

- $\tilde{\mathbf{x}}_t = [\mathbf{x}_t; \mathbf{s}]$
- $\tilde{\mathbf{x}}_t = \mathbf{x}_t + \mathbf{s}$
- Changes the input to the generation itself and not the context
- Not shown promising results so far

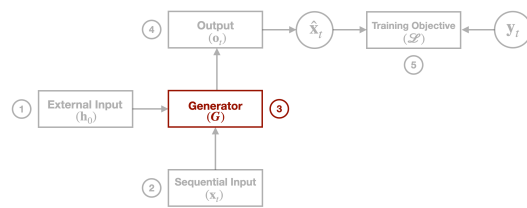


[Noraset et al. (2017), Zhou et al. (2018), Prabhumoye et al. (2019)]

Generator Operations

1. Controlled Generator Operations

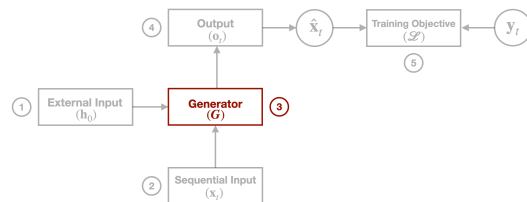
- $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t + \tanh(\mathbf{W}_d \mathbf{d}_t)$
 - \mathbf{d}_t = dialogue act representation, change made to LSTM cell
 - Add *dialogue act* information in the generation process
- $\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_h \mathbf{x}_t + \mathbf{r}_t \odot \mathbf{U}_h \mathbf{h}_{t-1} + \mathbf{s}_t \odot \mathbf{Yg} + \mathbf{q}_t \odot (\mathbf{1}_L^T \mathbf{Z} \mathbf{E}_t^{new})^T)$
 - \mathbf{s}_t = goal select gate; \mathbf{q}_t = item select gate, GRU cell
 - *recipe generation* task



[Gan et al. (2017), Kiddon et al. (2016), Wen et al. (2015)]

Generator Operations

1. Controlled Generator Operations



- $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t + \tanh(\mathbf{W}_d \mathbf{d}_t)$
 - \mathbf{d}_t = dialogue act representation, change made to LSTM cell
 - Add *dialogue act* information in the generation process
- $\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_h \mathbf{x}_t + \mathbf{r}_t \odot \mathbf{U}_h \mathbf{h}_{t-1} + \mathbf{s}_t \odot \mathbf{Yg} + \mathbf{q}_t \odot (\mathbf{1}_L^T \mathbf{Z} \mathbf{E}_t^{new})^T)$
 - \mathbf{s}_t = goal select gate; \mathbf{q}_t = item select gate, GRU cell
 - *recipe generation* task

[Gan et al. (2017), Kiddon et al. (2016), Wen et al. (2015)]

Generator Operations

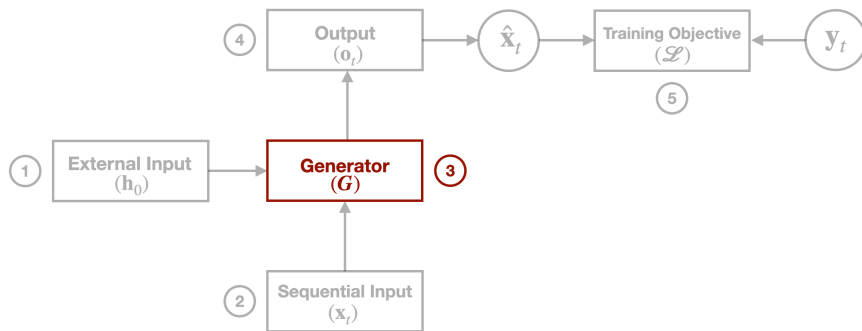
2. Recurrent Neural Networks

- LSTM, GRU

3. Transformers

4. Pre-trained language models

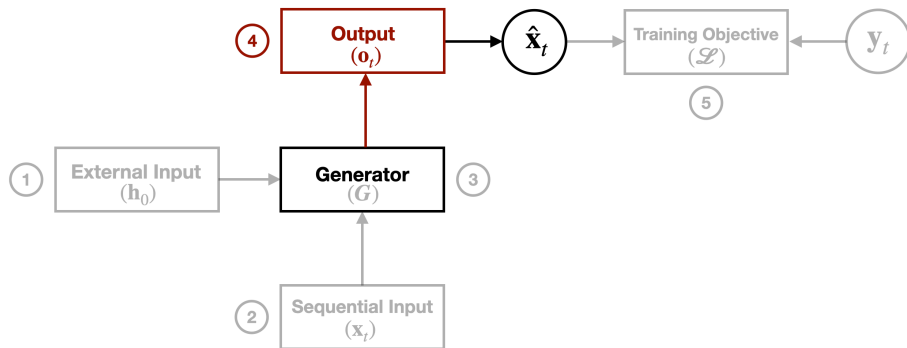
- BERT, GPT-2, BART, XL-Net



Output

1. Attention

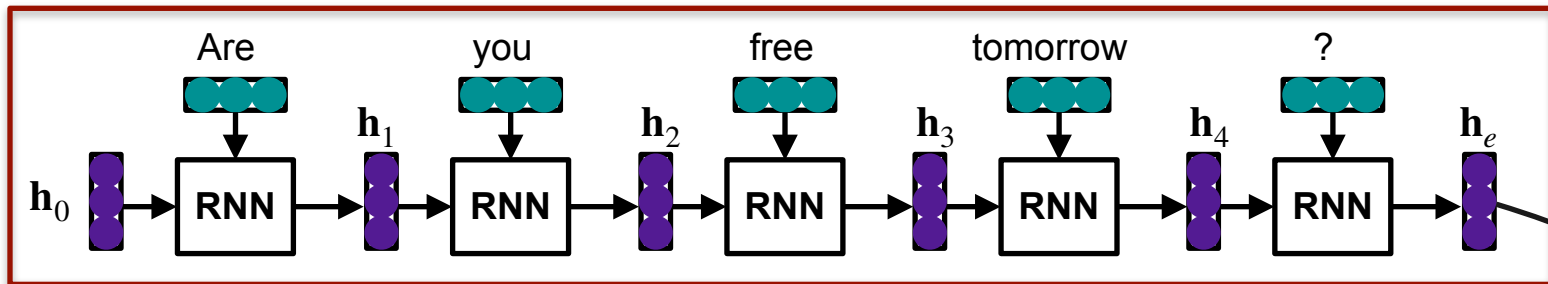
- Focus on source sequence
- Global Attention
- Local Attention
- Multi-headed Attention



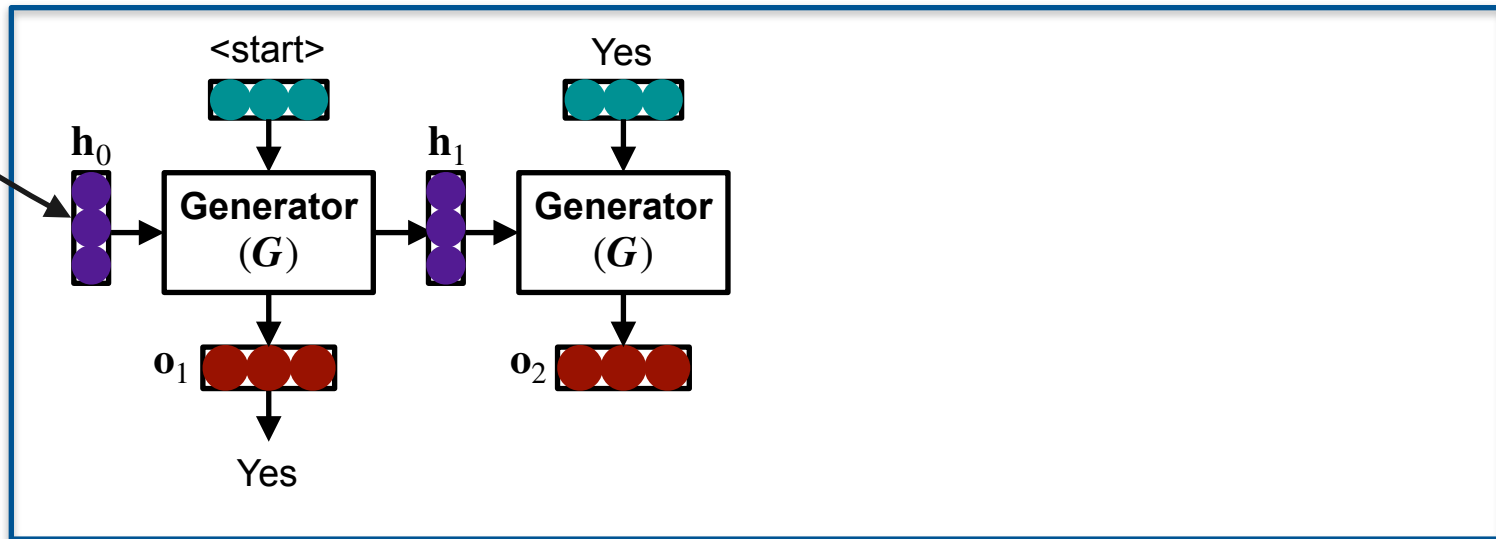
[Bahdanau et al. (2015), Luong et al. (2015), Vaswani et al. (2017)]

Generation Process

Encoder

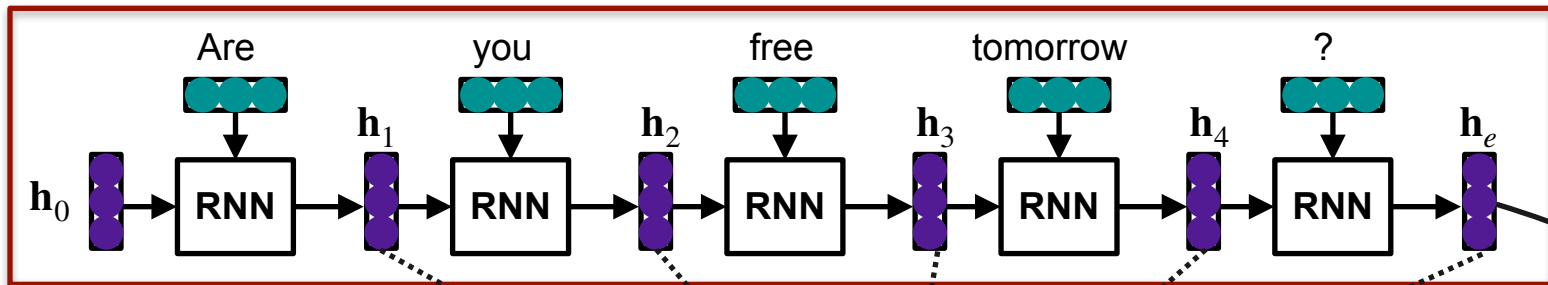


Decoder

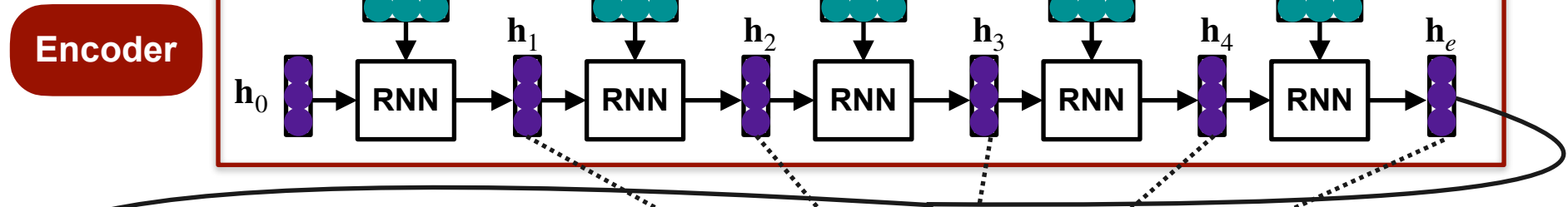
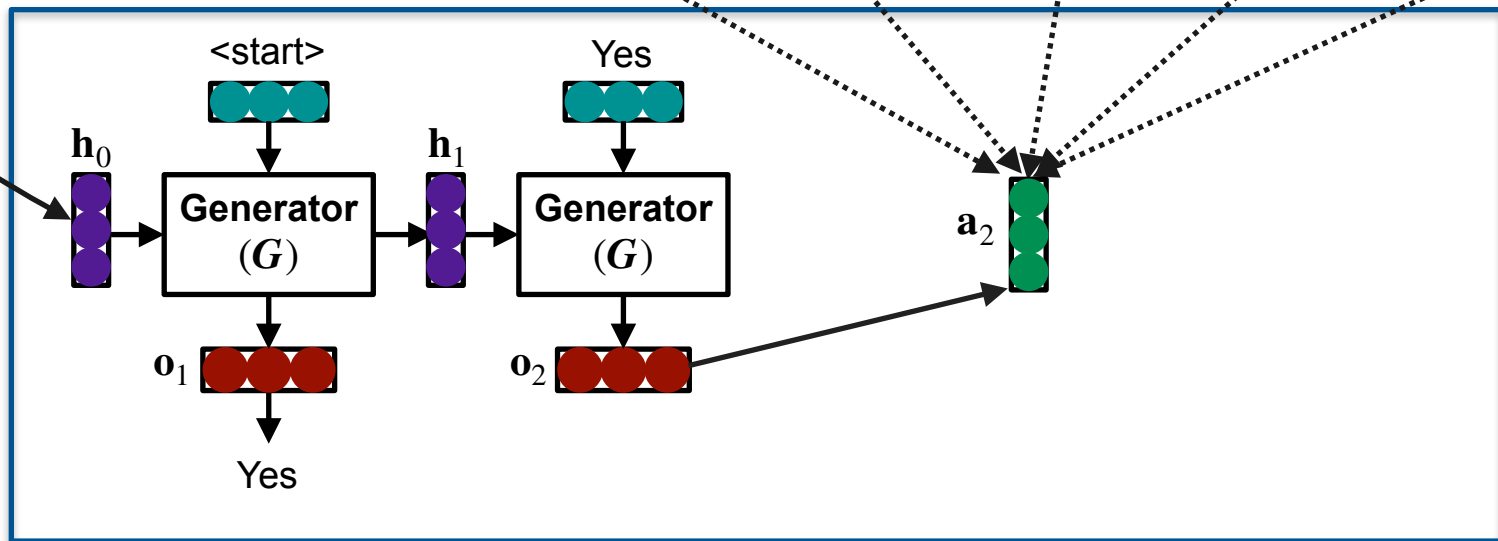


Generation Process

Encoder

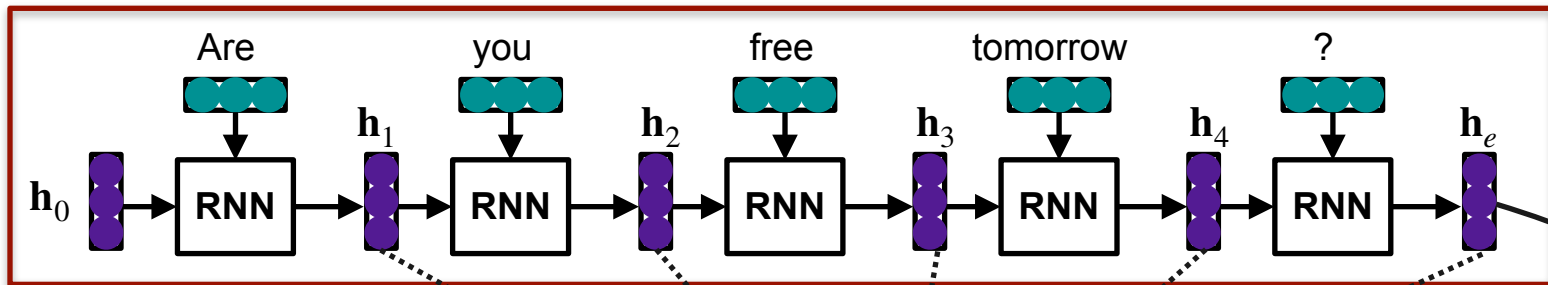


Decoder

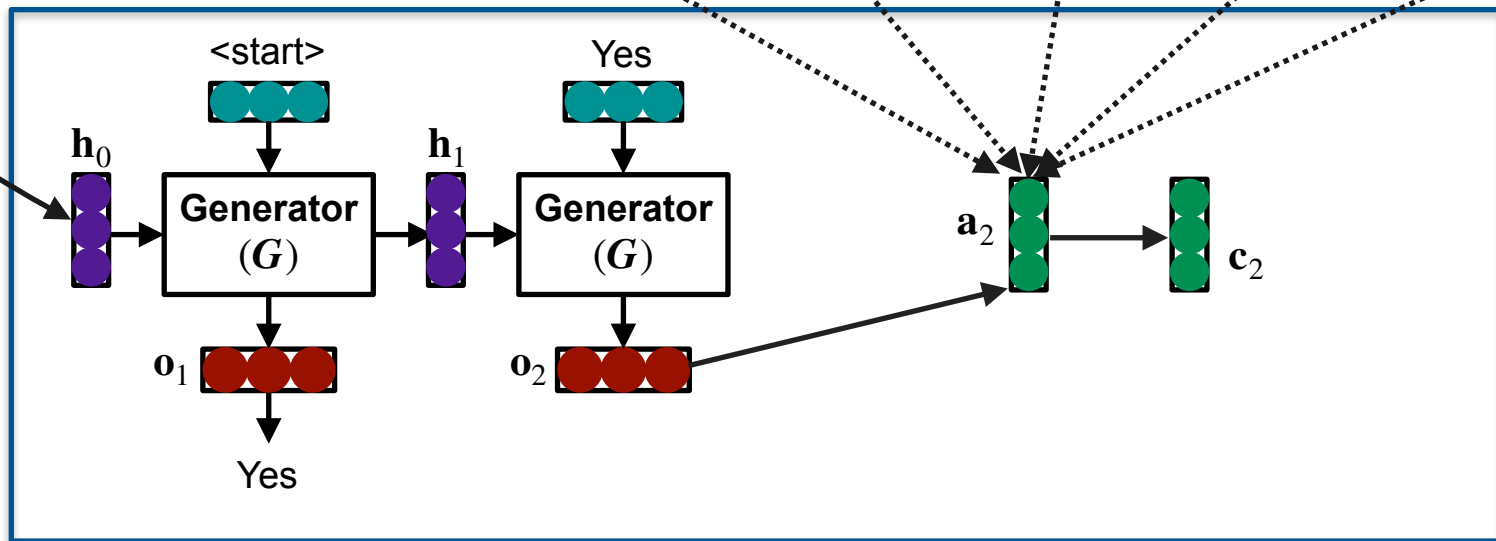


Generation Process

Encoder

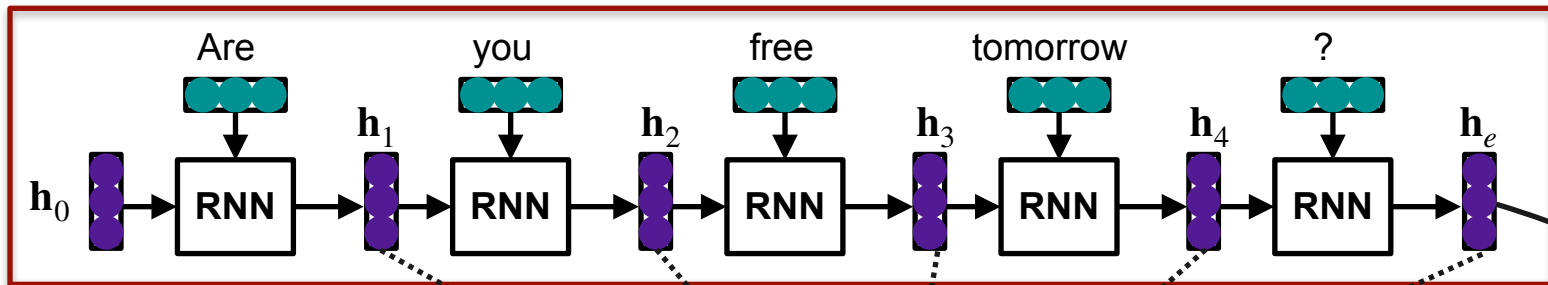


Decoder

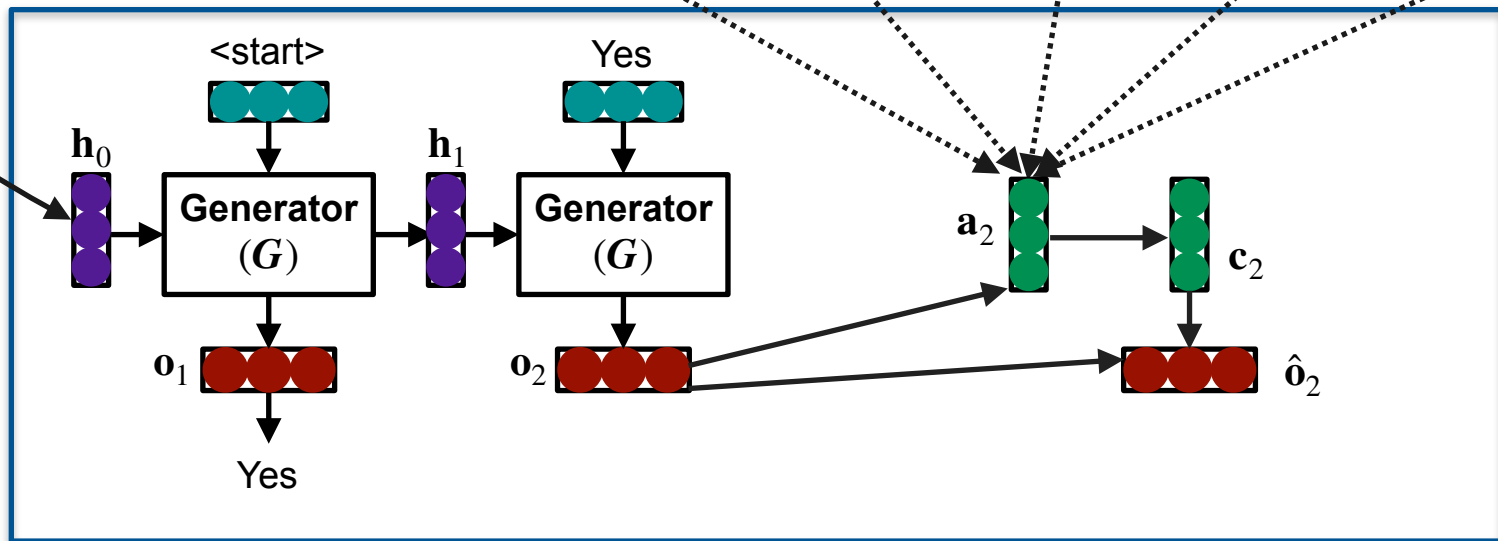


Generation Process

Encoder

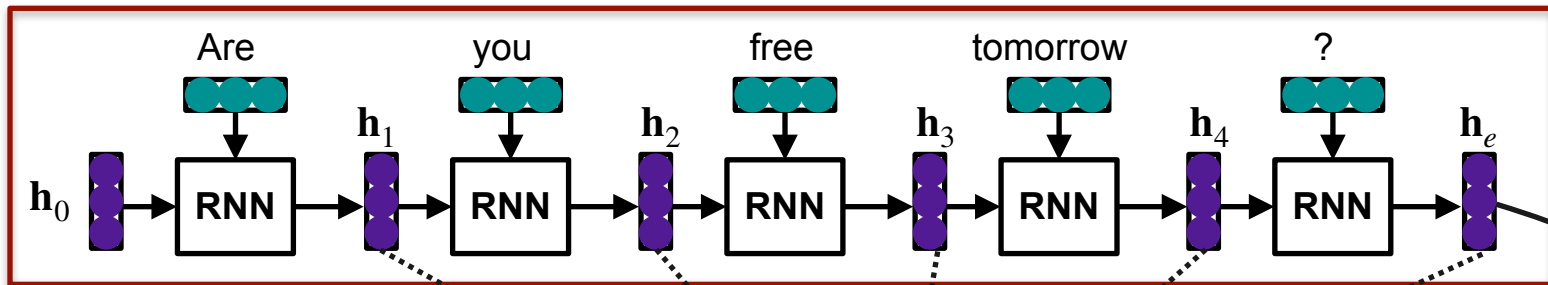


Decoder

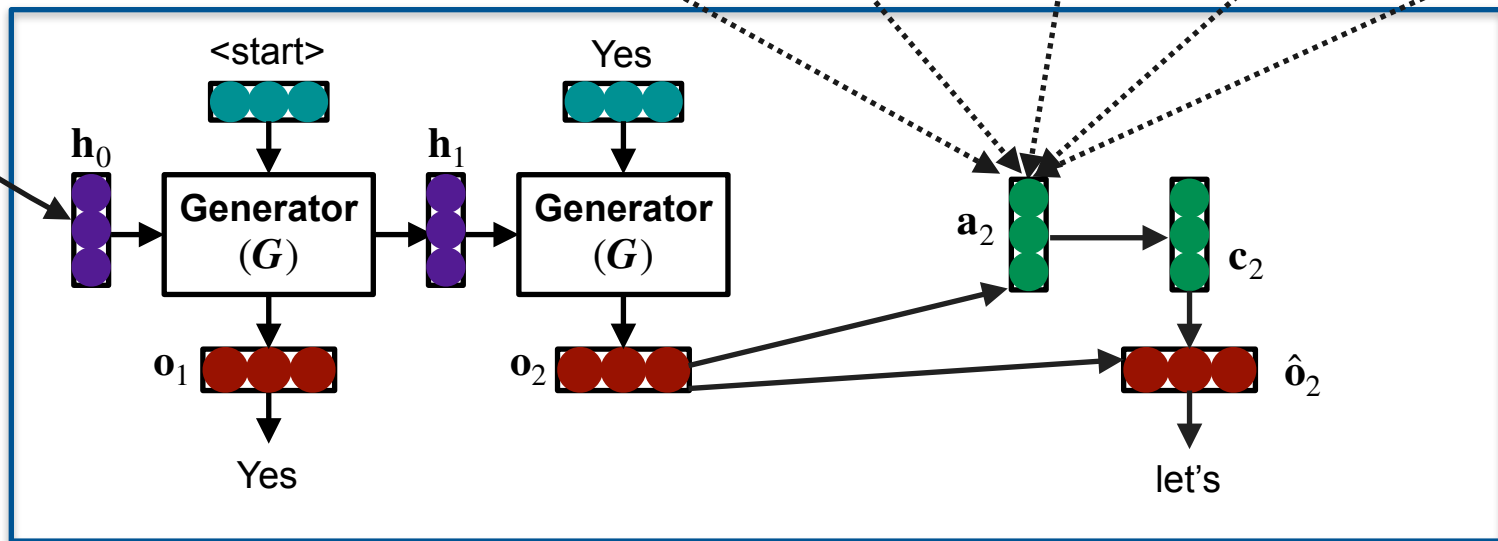


Generation Process

Encoder



Decoder



Output

1. Attention

- most effective - especially self and cross
- mostly control attribute tokens have been added to source sequence for attention
- under explored for controlling attributes but has a lot of potential

[Sudhakar et al. (2019), Dinan et al. (2018), Zhang et al. (2018)]

Output

2. External Feedback

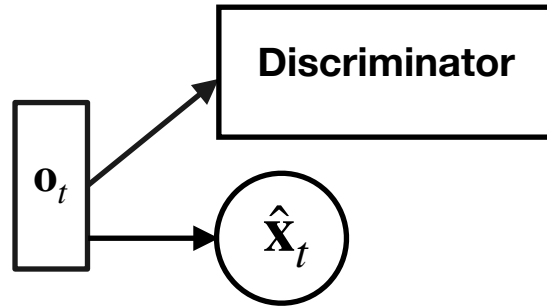
- discriminator has to be jointly trained like GAN

3. Arithmetic or linear transform

Output

2. External Feedback

- discriminator has to be jointly trained like GAN

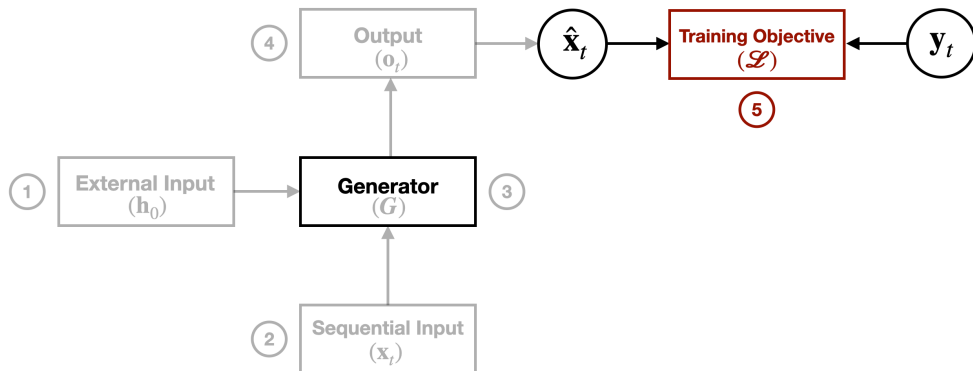


3. Arithmetic or linear transform

Training Objective

1. General Loss

- Cross Entropy Loss
- Unlikelihood Loss
- Decoding Strategies
- Used with any generation task



2. Classifier Loss

- design multiple classifier for any control attributes

[Welleck et al. (2020), Prabhumoye et al. (2018), Yang et al. (2018)]

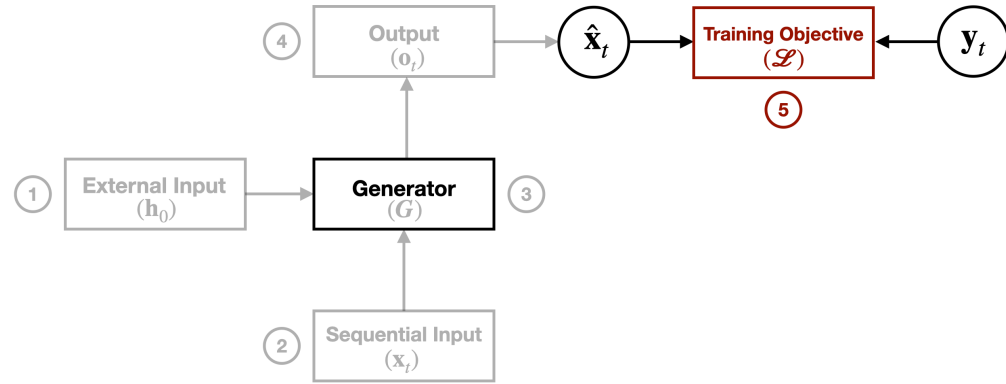
Training Objective

3. KL Divergence

- used with stochastic changes

4. Task Specific Loss

- design a loss for specific task (need not involve a classifier)
- Strategy Loss
- Coverage Loss
- Structure Loss



Future Work

- ***Empirical evaluation of schema***
 - to understand quantitatively which modules are more effective in controlling attributes
 - task-related architectures
 - add these control techniques to pre-trained models like BART, T5 etc