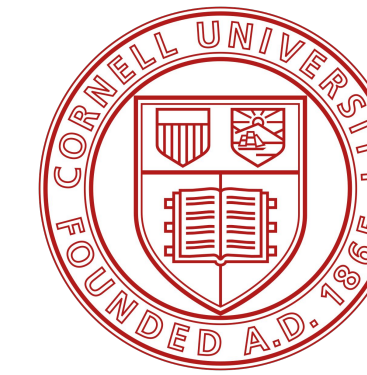# Block Diffusion: Interpolating between Autoregressive and Diffusion Language Models

Marianne Arriola[1], Aaron Gokaslan[1], Justin T. Chiu[2], Zhihan Yang[1], Zhixuan Qi[1], Jiaqi Han[3], Subham Sahoo[1], Volodymyr Kuleshov[1]

Cornell Tech[1], Cohere[2], Stanford University[3]

## 1 — AR or Diffusion for language? / Better together!

**Autoregression**

*Generation*

They're
They're speculating
They're speculating about . . .

❌ Sequential   ❌ Causal context only   ✅ High quality   ✅ Flexible-length   ✅ KV caching

**Diffusion**

*Generation*

Hirsh    need    account
Hirsh    need to    account data    to    released.
Hirsh will need to take account data that's to be released.

✅ Parallel   ✅ Global context   ❌ Low quality   ❌ Fixed-length   ❌ No KV caching

**Block Diffusion** *(Ours)*

*Generation*

Anatoly    well    for
Anatoly is well-known for his    plays, like
Anatoly is well-known for his witty chess plays, like    bold    Gambit . . .

✅ Parallel   ✅ Semi-global context   ✅ High quality   ✅ Flexible-length   ✅ KV caching

## 2 — Unifying AR and Diffusion Objectives

$$-\log p_\theta(\mathbf{x}) = -\sum_{b=1}^{B} \log p_\theta(\mathbf{x}^b|\mathbf{x}^{<b}) \leq \sum_{b=1}^{B} \mathcal{L}_{\text{diff}}(\mathbf{x}^b|\mathbf{x}^{<b};\theta)$$

Clean context $\mathbf{x}^{<b}$

Anatoly is well-known for

Clean block $\mathbf{x}^b$

his witty chess plays, like

**BERT-style training within each block**

1) Mask with probability $t \sim [0,1]$

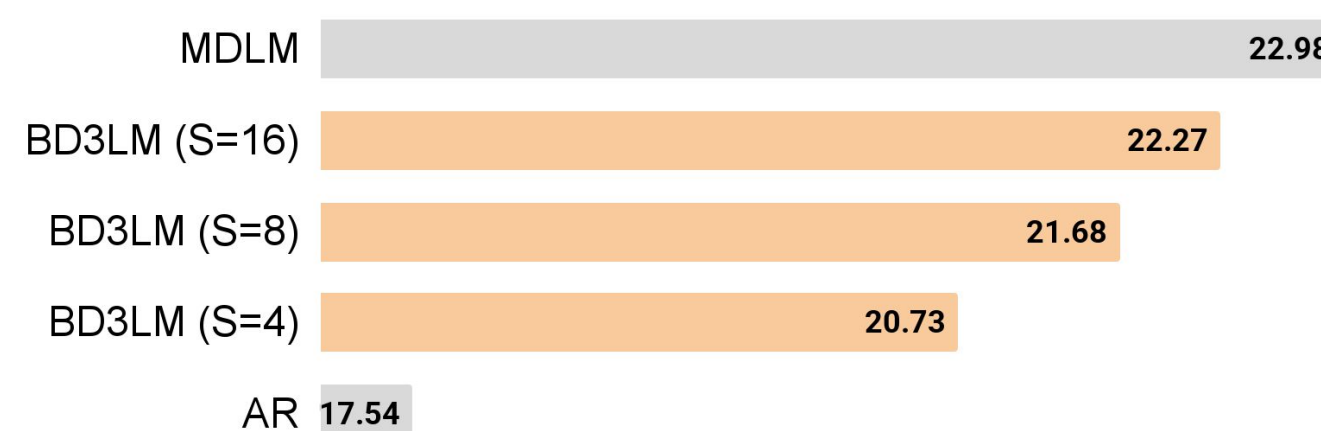2) Unmask $p_\theta(\mathbf{x}^b|\mathbf{x}_t^b, \mathbf{x}^{<b})$

**Diffusion loss $\mathcal{L}_{\text{diff}}$**

$$\mathbb{E}_{t,\mathbf{x}_t \sim q} -\frac{1}{t}\log p_\theta(\mathbf{x}^b|\mathbf{x}_t^b, \mathbf{x}^{<b})$$
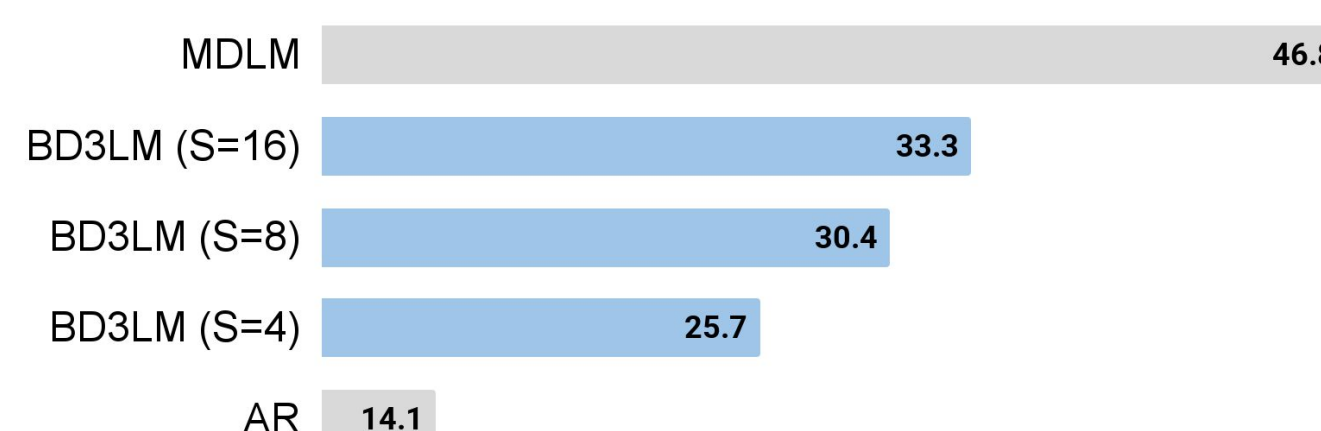
his ◆ chess ◆ like

Noised block $\mathbf{x}_t^b$

## 3 — Predict all blocks in 1 pass

**Goal:** Predict $B$ blocks

Predictions

*Block 1*    Transformer

Inputs    ◆ World!

*Block 2*    Transformer

Hello World!    Today ◆

**Solution:** Use 1 pass with a custom attention pattern

Transformer

Noised blocks $\mathbf{x}_{t_1}^1 \oplus \cdots \oplus \mathbf{x}_{t_B}^B$
◆ World! Today ◆ ◆ to

Clean blocks $\mathbf{x}^{1:B}$
Hello World! Today I'm going to

**Attention Pattern**

## 4 — Arbitrary-Length Generation + KV caching

*Block 1*

Samples    At ◆ the ◆
**Sampler**
Inputs    ◆ ◆ ◆ ◆

x Generation steps

*Block 2*

told the ◆ that
**Sampler**
KV cache    ◆ ◆ ◆ ◆

x Generation steps . . .

Accumulated Samples

At 9AM, the coach

## 5 — SOTA Likelihoods & Long Document Gen.

**S = Block size**

**Test Perplexity (↓) on OpenWebText**

| | |
|---|---|
| MDLM | 22.98 |
| BD3LM (S=16) | 22.27 |
| BD3LM (S=8) | 21.68 |
| BD3LM (S=4) | 20.73 |
| AR | 17.54 |

**Perplexity (↓) of Generated Samples under GPT2**

| | |
|---|---|
| MDLM | 46.8 |
| BD3LM (S=16) | 33.3 |
| BD3LM (S=8) | 30.4 |
| BD3LM (S=4) | 25.7 |
| AR | 14.1 |

**Maximum Document Generation Length (↑)**

| | Length (tokens) |
|---|---|
| SEDD | 1024 |
| BD3-LM (S=16) | 9982 |