

# Обучение рекуррентных нейронных сетей с использованием сопряженных уравнений

**Студент:** Хоменко Р.Д.

**Руководитель:** Ревизников Д.Л., профессор-совм. кафедры 810Б, д.ф.-м.н., профессор

# Постановка задачи

Исследовать возможность применения Neural ODEs в задачах обработки естественного языка. Сравнить данный подход с классическими архитектурами на примере модели генерации текста.

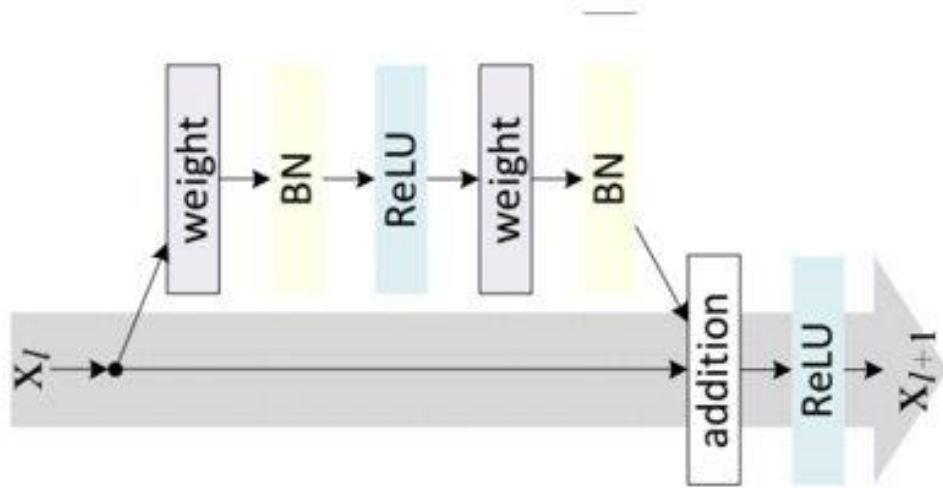
## Обзор исследований

- Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, David Duvenaud. **Neural Ordinary Differential Equations**
- Hanshu Yan, Jiawei Du, Vincent Y. F. Tan, Jiashi Feng. **On Robustness of Neural Ordinary Differential Equations**
- Edward De Brouwer, Jaak Simm, Adam Arany, Yves Moreau. **GRU-ODE-Bayes: Continuous modeling of sporadically-observed time series**
- Prasad Kawthekar, Raunaq Rewari, Suvrat Bhooshan. **Evaluating Generative Models for Text Generation**

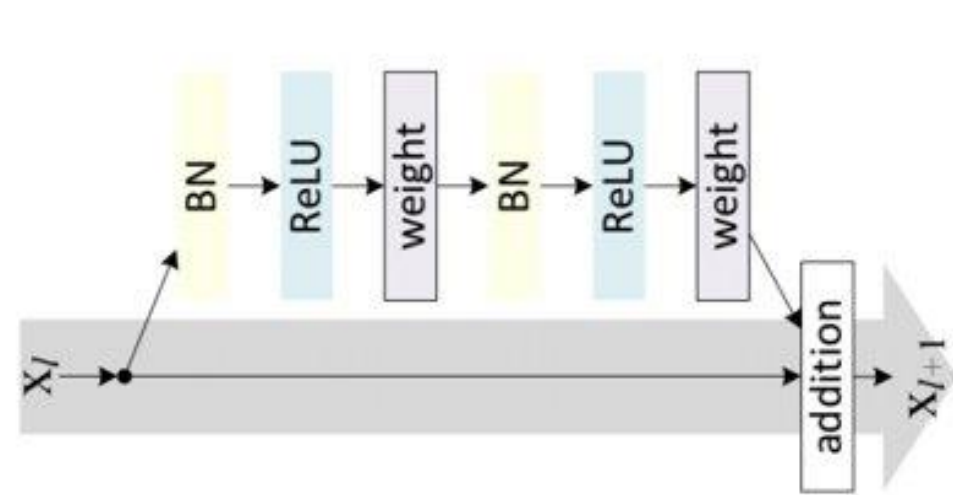
## Актуальность темы исследования

- Ускорение обучения на больших объемах данных
- Контроль над точностью и вычислительными ресурсами
- Исследование комбинации Neural ODE с другими слоями
- Исследование модели (устойчивость решения и т.д.)

# Архитектура ResNet

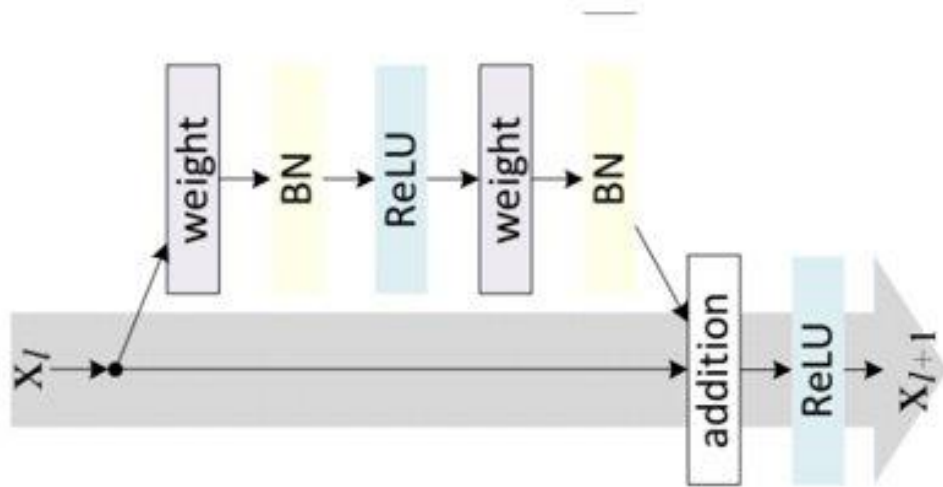


(a) ResNet блок

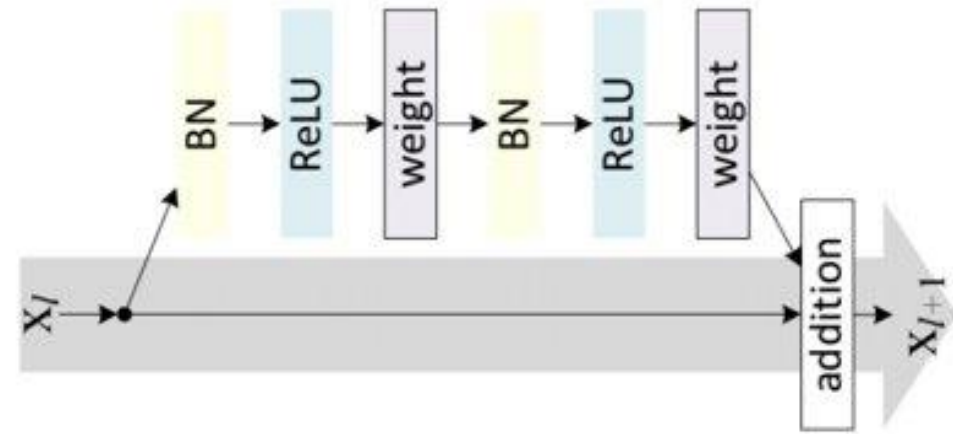


(б) ResNet блок с пред-активацией

# Архитектура ResNet



(a) ResNet блок



(б) ResNet блок с пред-активацией

$$\begin{cases} \frac{dx}{dt} = f(x, t), \\ x(0) = x_0 \end{cases}$$

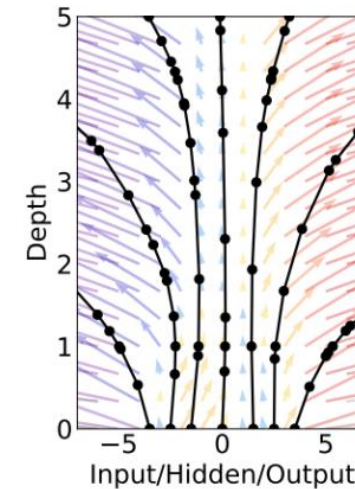
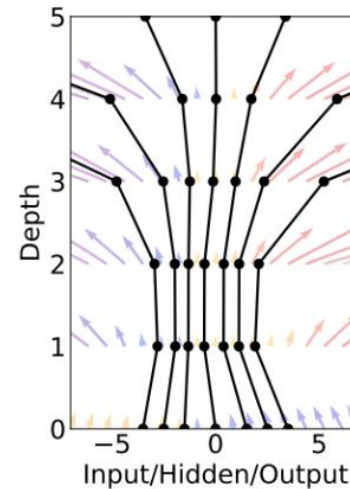
$$x_{l+1} = x_l + hf(x_l, t_l)$$

# Neural ODEs

$$\begin{cases} \frac{dz}{dt} = f(z(t), \theta), & t \in [t_0, t_1] \\ z(t_0) = z_0 \end{cases}$$

$$L(z(t_1)) \rightarrow \min_{\theta}$$

- $z_0$  - вектор признаков
- $z(t_1)$  - выход сети
- $\theta$  – вектор параметров модели
- $L$  – функция потерь
- $z(t) = z(t_0) + \int_{t_0}^t f(z(t), t, \theta) dt$



# Аjoint метод: условия оптимальности (1)

$$L(z(t_1), \theta) \quad \frac{dz}{dt} = f(z(t), \theta), \quad z(t_0) = z_0$$

$$\mathcal{L}(z(t), \theta, a(t)) = L(z(t_1), \theta) + \int_{t_0}^{t_1} a(t)^\top \left( \frac{dz}{dt} - f(z(t), \theta) \right) dt + b \cdot (z(t_0) - z_0)$$

$\theta$ :

$$\delta_\theta \mathcal{L}(\theta, h) = \frac{\partial L}{\partial \theta}(z(t_1), \theta) \cdot h - \int_{t_0}^{t_1} a(t)^\top \frac{\partial f}{\partial \theta}(z(t), t, \theta) dt \cdot h = 0 \quad \text{for any feasible variation } h$$

$$\begin{cases} \frac{d}{dt} \left( \frac{\partial L}{\partial \theta} \right) = -a(t)^\top \frac{\partial f}{\partial \theta}(z(t), \theta) \\ \left. \frac{\partial L}{\partial \theta} \right|_{t=t_1} = 0 \end{cases}$$

$$\delta J(y, h) \triangleq \lim_{\varepsilon \rightarrow 0} \frac{J[y + \varepsilon h] - J[y]}{\varepsilon} = \left. \frac{\partial}{\partial \varepsilon} J[y + \varepsilon h] \right|_{\varepsilon=0}$$



## Аjoint метод: условия оптимальности (2)

$$\mathcal{L}(z(t), \theta, a(t)) = L(z(t_1), \theta) + \int_{t_0}^{t_1} a(t)^\top \left( \frac{dz}{dt} - f(z(t), \theta) \right) dt + b \cdot (z(t_0) - z_0)$$

z:

$$\delta_z \mathcal{L}(z, h) = \frac{\partial L}{\partial z}(z(t_1), \theta) \cdot h(t_1) + \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \int_{t_0}^{t_1} a(t)^\top \left( \varepsilon \frac{dh(t)}{dt} - \frac{\partial f}{\partial z}(z(t), \theta) \cdot \varepsilon h(t) + o(\varepsilon) \right) dt$$

$$\int_{t_0}^{t_1} a(t)^\top \frac{dh(t)}{dt} dt = \left[ a(t)^\top h(t) \right] \Big|_{t_0}^{t_1} - \int_{t_0}^{t_1} \frac{da(t)}{dt} h(t) dt = -a(t_1)^\top h(t_1) - \int_{t_0}^{t_1} \frac{da(t)}{dt} h(t) dt$$

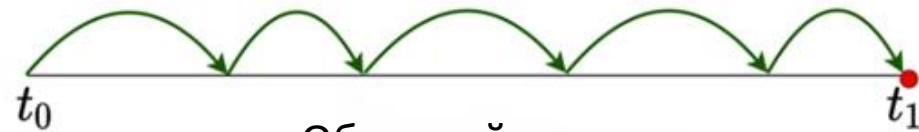
$$\delta_z \mathcal{L}(z, h) = \left( \frac{\partial L}{\partial z}(z(t_1)) - a(t_1) \right) \cdot h(t_1) - \int_{t_0}^{t_1} \left\{ \frac{da(t)}{dt} + a(t)^\top \frac{\partial f}{\partial z}(z(t), \theta) \right\} h(t) dt = 0$$

$$\frac{da(t)}{dt} = -a(t)^\top \frac{\partial f}{\partial z}(z(t), \theta), \quad a(t_1) = \frac{\partial L}{\partial z}(z(t_1))$$

# Аjoint метод

$$\left\{ \begin{array}{l} \frac{d}{dt} \frac{\partial L}{\partial \theta} = -a(t)^\top \frac{\partial f}{\partial \theta}(z(t), \theta) \\ \frac{\partial L}{\partial \theta} \Big|_{t=t_1} = \mathbf{0} \end{array} \right. \quad \boxed{\left\{ \begin{array}{l} \frac{dz}{dt} = f(z(t), \theta) \\ z(t_1) = z_1 \end{array} \right.} \quad \left\{ \begin{array}{l} \frac{da(t)}{dt} = -a(t)^\top \frac{\partial f}{\partial z}(z(t), \theta) \\ a(t_1) = \frac{\partial L}{\partial z}(z(t_1)) \end{array} \right.$$

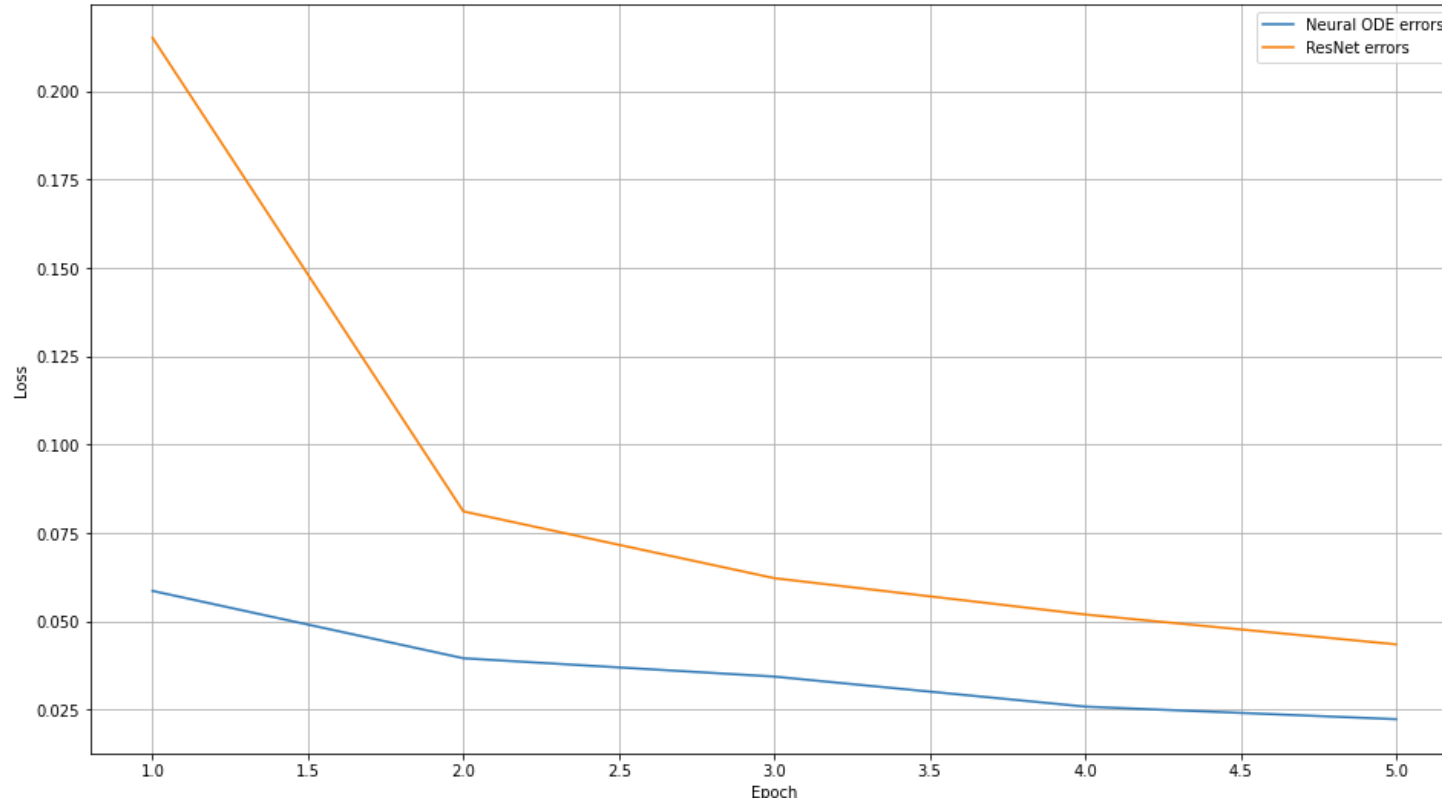
Прямой проход



Обратный проход



# Аjoint метод: сравнение на архитектуре ResNet



# GRU-ODE

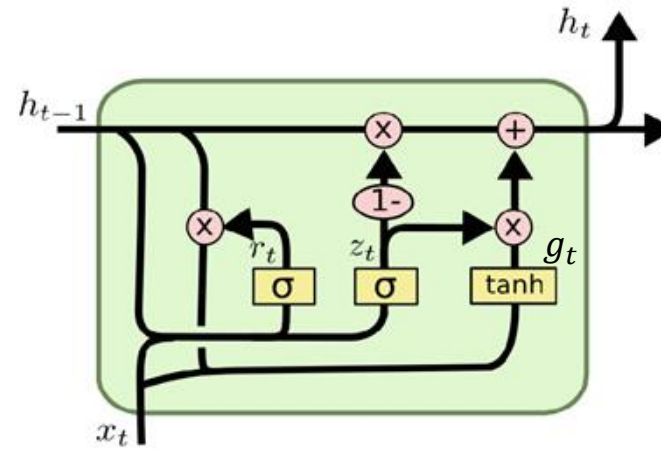
## GRU

$$\mathbf{h}_t = \mathbf{z}_t \odot \mathbf{h}_{t-1} + (1 - \mathbf{z}_t) \odot \mathbf{g}_t$$

$$\mathbf{r}_t = \sigma(W_r \mathbf{x}_t + U_r \mathbf{h}_{t-1} + \mathbf{b}_r)$$

$$\mathbf{z}_t = \sigma(W_z \mathbf{x}_t + U_z \mathbf{h}_{t-1} + \mathbf{b}_z)$$

$$\mathbf{g}_t = \tanh(W_h \mathbf{x}_t + U_h (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_h)$$



## GRU-ODE

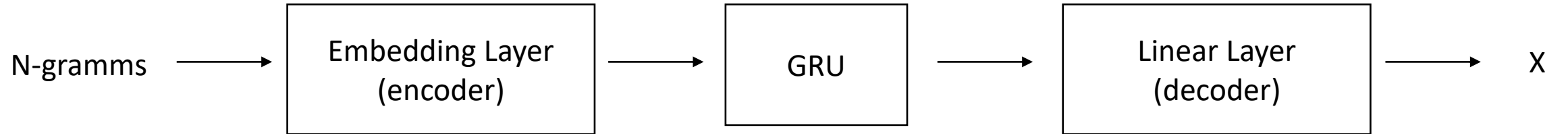
$$\frac{d\mathbf{h}(t)}{dt} = (1 - \mathbf{z}(t)) \odot (\mathbf{g}(t) - \mathbf{h}(t))$$

$$\mathbf{r}(t) = \sigma(W_r \mathbf{x}(t) + U_r \mathbf{h}(t) + \mathbf{b}_r)$$

$$\mathbf{z}(t) = \sigma(W_z \mathbf{x}(t) + U_z \mathbf{h}(t) + \mathbf{b}_z)$$

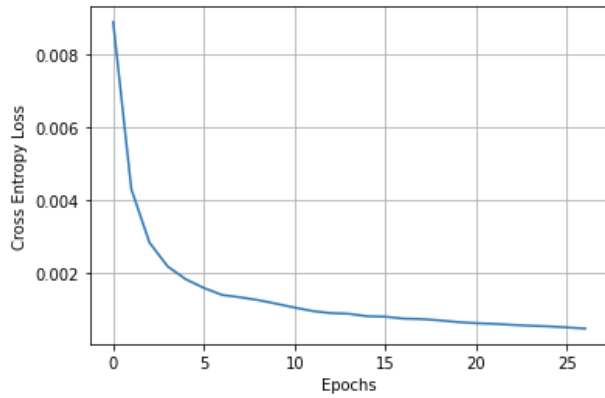
$$\mathbf{g}(t) = \tanh(W_h \mathbf{x}(t) + U_h (\mathbf{r}(t) \odot \mathbf{h}(t)) + \mathbf{b}_h)$$

# Модель генерации текста

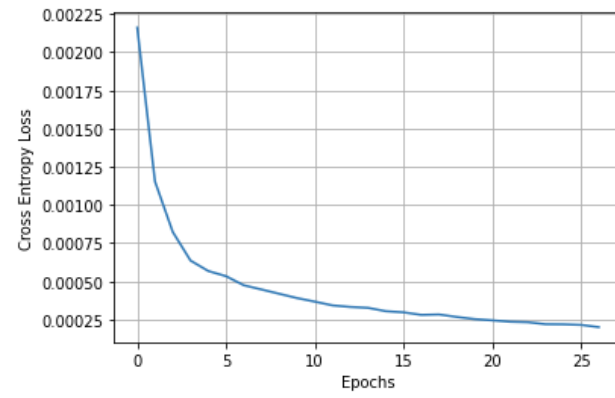


$$i_{word} = \text{multinomial}(\exp\left(\frac{X}{T}\right), 1)$$

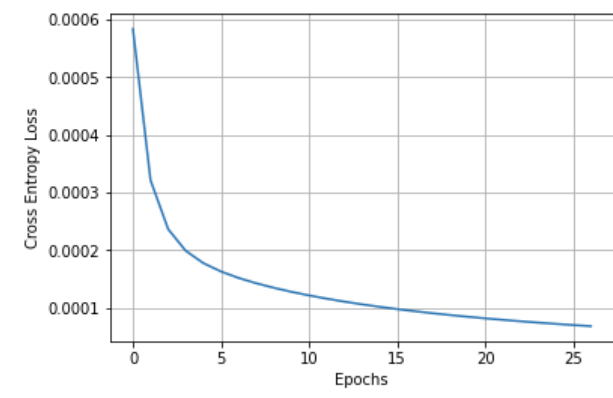
# Обучение модели



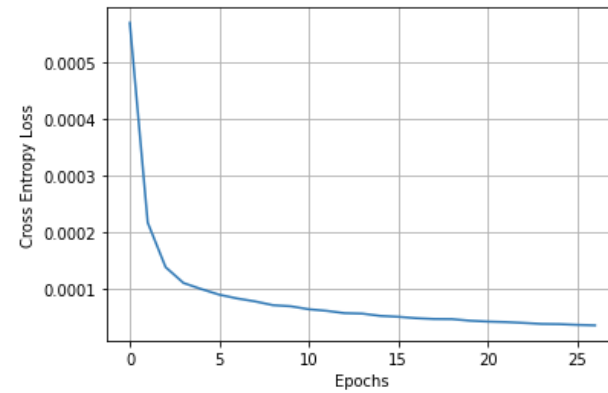
vanilla RNN



GRU

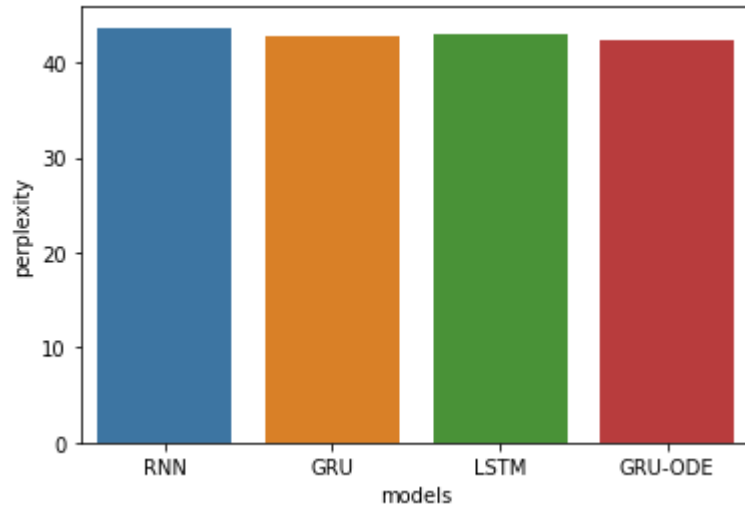


LSTM

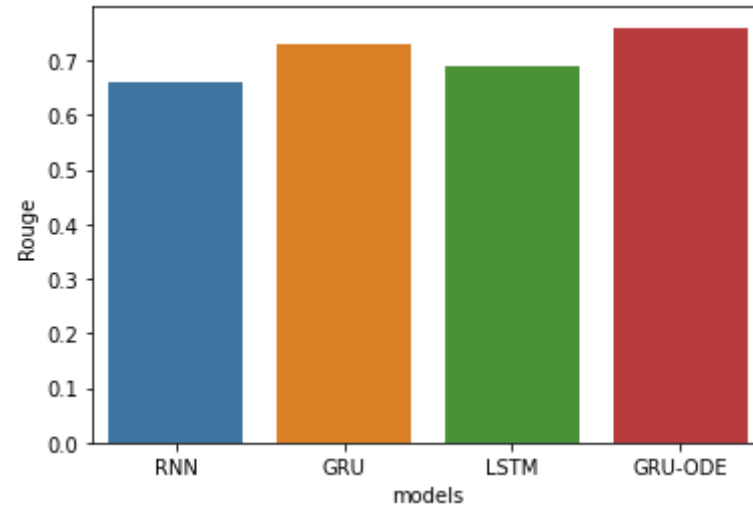


GRU-ODE

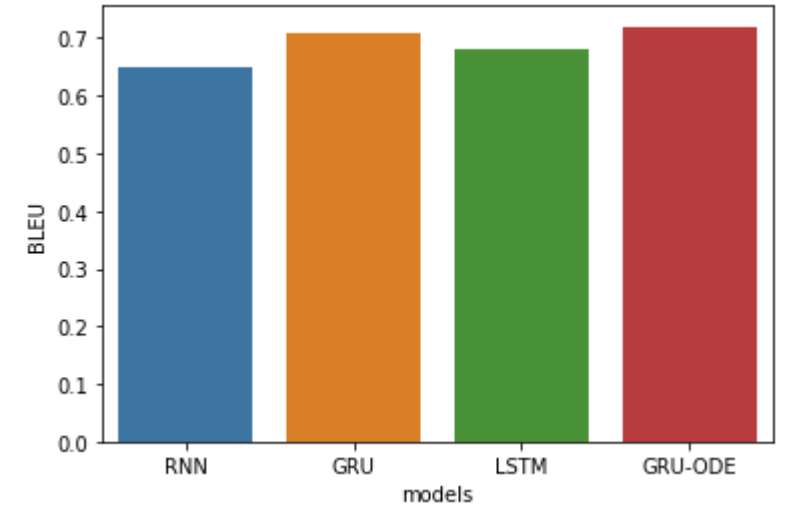
# Метрики качества



Perplexity



Rouge



BLEU

# Проблемы

- «жесткость» систем при backward pass
- некорректные результаты при использовании неявных методов
- оценка качества генеративных моделей



# Методы решения

- применение специфических методов для backward pass
- использование методов с адаптивным шагом
- комбинированные метрики

## Заключение

- Исследована модель нейронных ОДУ
- Реализован алгоритм обучения GRU-ячейки на основе нейронных ОДУ
- Реализация успешно протестирована на модели генерации текста

Нейронные ОДУ могут быть успешно интегрированы в существующих моделях.