

# Introduction to Intelligent Systems / Lecture 4

---

## Feedback Neural Networks (a.k.a. Recurrent neural networks)

Jae Yun JUN KIM\*

July 9, 2019

A recurrent neural network (RNN) is a class of artificial neural network where connections between units form a directed cycle.

This creates an internal state of the network which allows it to exhibit dynamic temporal behavior.

Unlike feedforward neural networks, RNNs can use their internal memory to process arbitrary sequences of inputs.

Possible applications: handwriting recognition / speech recognition.

**Example:** XOR operation.

Consider its truth table:

Input		Output
0	0	0
0	1	1
1	0	1
1	1	0

Table 1: The truth table for the XOR operation

Then, let us create a sequence from the above table

000011101110.

If we did not know where this sequence came from, it would have been very difficult to find its structure. From this observation, a question arises naturally: how to learn automatically the structure of a given sequence of data, so that we can make predictions?

We might be able to answer this question using an RNN.

## 1 Simple recurrent neural networks

There are two well-known types of RNNs: Elman RNN and Jordan RNN.

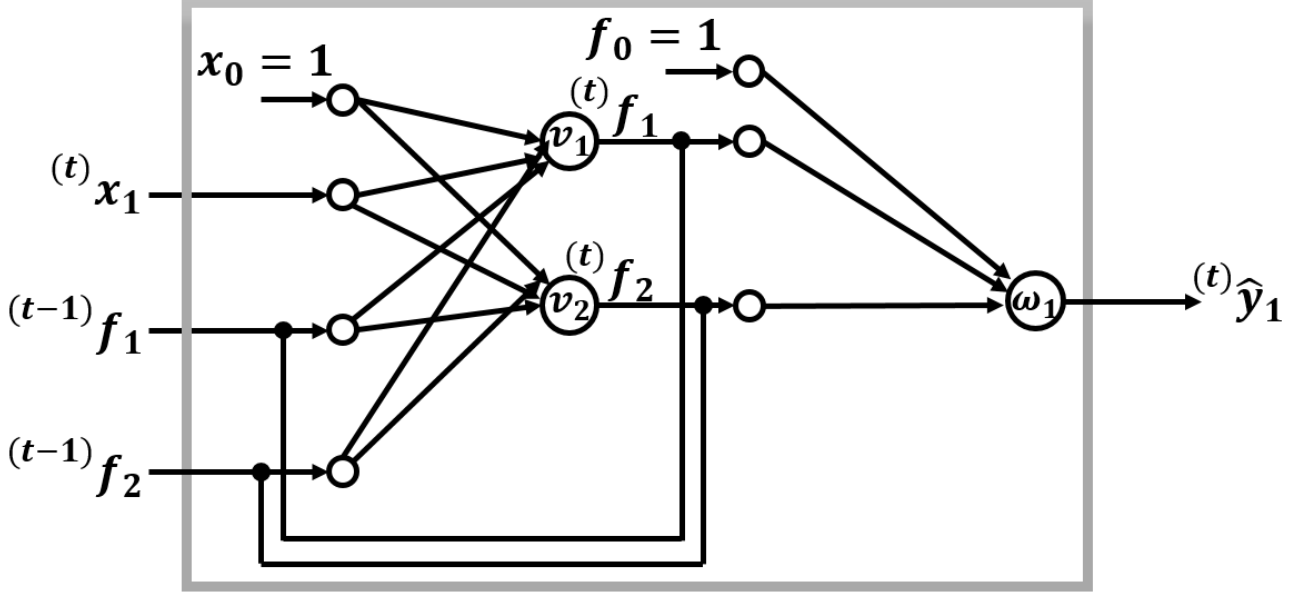


Figure 1: Elman RNN

## 1.1 Elman RNN

### 1.1.1 Forward propagation

Input data matrix

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(I)})^T \end{bmatrix} \in \mathbb{R}^{I \times N}, \quad \text{where } x^{(i)} \in \mathbb{R}^{N \times 1} \quad (1)$$

Extended input data matrix

$$\bar{X} = [\mathcal{I} \quad X \quad (t-1)F] = \begin{bmatrix} (\bar{x}^{(1)})^T \\ \vdots \\ (\bar{x}^{(I)})^T \end{bmatrix} \in \mathbb{R}^{I \times (N+1+K)}, \quad \text{where } \mathcal{I} = [1, \dots, 1]^T \in \mathbb{R}^{I \times 1} \quad (2)$$

$$\bar{\bar{X}} = \bar{X} \cdot v = \begin{bmatrix} (\bar{x}^{(1)})^T \cdot v \\ \vdots \\ (\bar{x}^{(I)})^T \cdot v \end{bmatrix} \in \mathbb{R}^{I \times K} \quad (3)$$

Define

$$F = \left(1 + \exp\left(-\bar{\bar{X}}\right)\right)^{-1} = \begin{bmatrix} (f^{(1)})^T \\ \vdots \\ (f^{(I)})^T \end{bmatrix} \in \mathbb{R}^{I \times K} \quad (4)$$

---

\*ECE Paris Graduate School of Engineering, 37 quai de Grenelle 75015 Paris, France; jae-yun.jun-kim@ece.fr

$$\bar{F} = [\mathcal{I} \quad F] = \begin{bmatrix} (\bar{f}^{(1)})^T \\ \vdots \\ (\bar{f}^{(I)})^T \end{bmatrix} \in \mathbb{R}^{I \times (K+1)}, \quad \text{where } \mathcal{I} = [1, \dots, 1]^T \in \mathbb{R}^{I \times 1} \quad (5)$$

$$\bar{\bar{F}} = \bar{F} \cdot \omega = \begin{bmatrix} (\bar{f}^{(1)})^T \cdot \omega \\ \vdots \\ (\bar{f}^{(I)})^T \cdot \omega \end{bmatrix} \in \mathbb{R}^{I \times J} \quad (6)$$

Define

$$G = \left(1 + \exp\left(-\bar{\bar{F}}\right)\right)^{-1} = \begin{bmatrix} (g^{(1)})^T \\ \vdots \\ (g^{(I)})^T \end{bmatrix} \in \mathbb{R}^{I \times J} \quad (7)$$

## 1.2 Jordan RNN

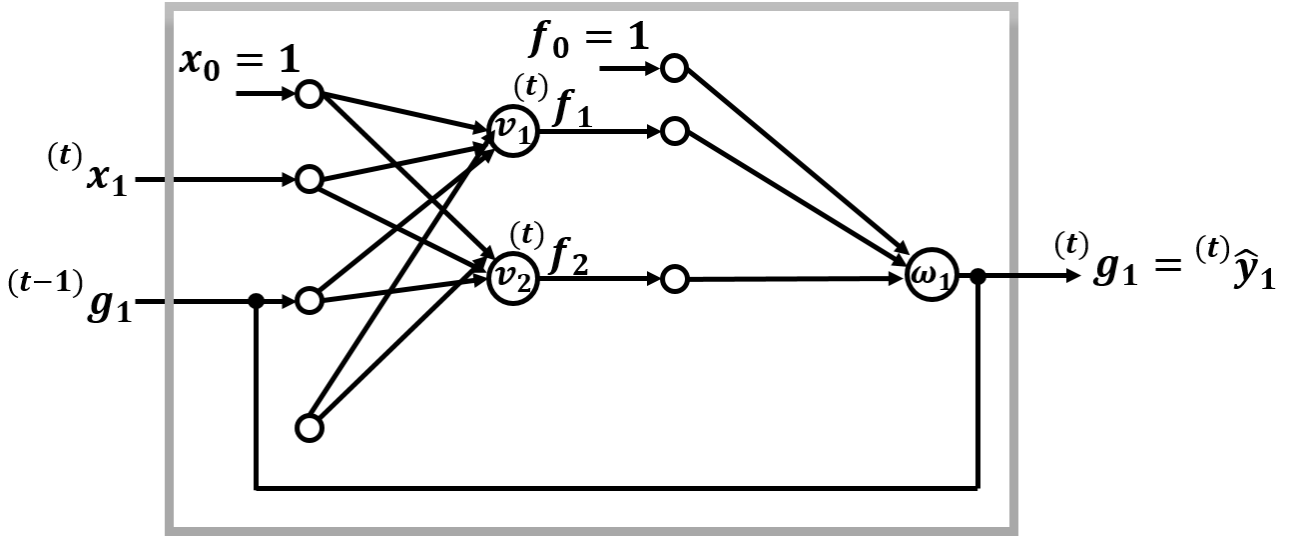


Figure 2: Jordan RNN

### 1.2.1 Forward propagation

Input data matrix

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(I)})^T \end{bmatrix} \in \mathbb{R}^{I \times N}, \quad \text{where } x^{(i)} \in \mathbb{R}^{N \times 1} \quad (8)$$

Extended input data matrix

$$\overline{X} = [\mathcal{I} \quad X \quad {}^{(t-1)}G] = \begin{bmatrix} (\overline{x}^{(1)})^T \\ \vdots \\ (\overline{x}^{(I)})^T \end{bmatrix} \in \mathbb{R}^{I \times (N+1+J)}, \quad \text{where } \mathcal{I} = [1, \dots, 1]^T \in \mathbb{R}^{I \times 1} \quad (9)$$

$$\overline{\overline{X}} = \overline{X} \cdot v = \begin{bmatrix} (\overline{x}^{(1)})^T \cdot v \\ \vdots \\ (\overline{x}^{(I)})^T \cdot v \end{bmatrix} \in \mathbb{R}^{I \times K} \quad (10)$$

Define

$$F = \left(1 + \exp\left(-\overline{\overline{X}}\right)\right)^{-1} = \begin{bmatrix} (f^{(1)})^T \\ \vdots \\ (f^{(I)})^T \end{bmatrix} \in \mathbb{R}^{I \times K} \quad (11)$$

$$\overline{F} = [\mathcal{I} \quad F] = \begin{bmatrix} (\overline{f}^{(1)})^T \\ \vdots \\ (\overline{f}^{(I)})^T \end{bmatrix} \in \mathbb{R}^{I \times (K+1)}, \quad \text{where } \mathcal{I} = [1, \dots, 1]^T \in \mathbb{R}^{I \times 1} \quad (12)$$

$$\overline{\overline{F}} = \overline{F} \cdot \omega = \begin{bmatrix} (\overline{f}^{(1)})^T \cdot \omega \\ \vdots \\ (\overline{f}^{(I)})^T \cdot \omega \end{bmatrix} \in \mathbb{R}^{I \times J} \quad (13)$$

Define

$$G = \left(1 + \exp\left(-\overline{\overline{F}}\right)\right)^{-1} = \begin{bmatrix} (g^{(1)})^T \\ \vdots \\ (g^{(I)})^T \end{bmatrix} \in \mathbb{R}^{I \times J} \quad (14)$$

## 2 Example

Suppose that we have a sequence of eight bits:  $XXXXXXX$ , where  $X \in \{0, 1\}$ . Then,

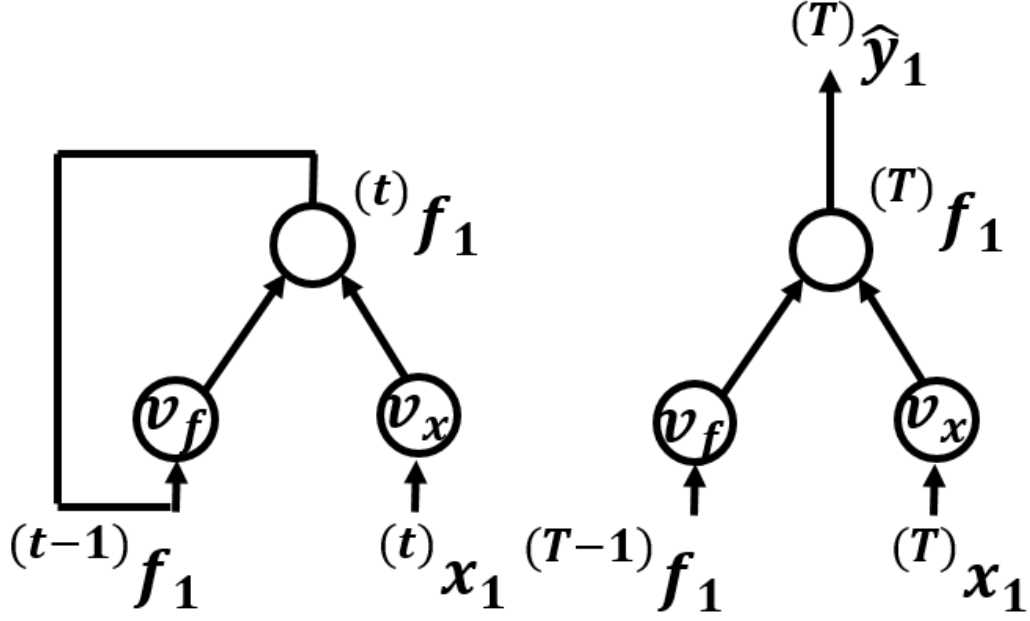
$$\hat{y} = \sum_{t=1}^T x_t \quad (15)$$

where and  $T = 8$ .

**Note:**

Given the fact that

- If  $\forall t, x_t = 0$ , then  $\hat{y} = 0$
- The desired output is linear with respect to the input, the activation function is designed to be the identity function.



(a) For intermediate states: for  $t = 1 : 7$       (b) For the last state: for  $t = T = 8$

Figure 3: Elman RNN example

- Consider the Elman RNN for simplicity

As a result, we obtain the following graph

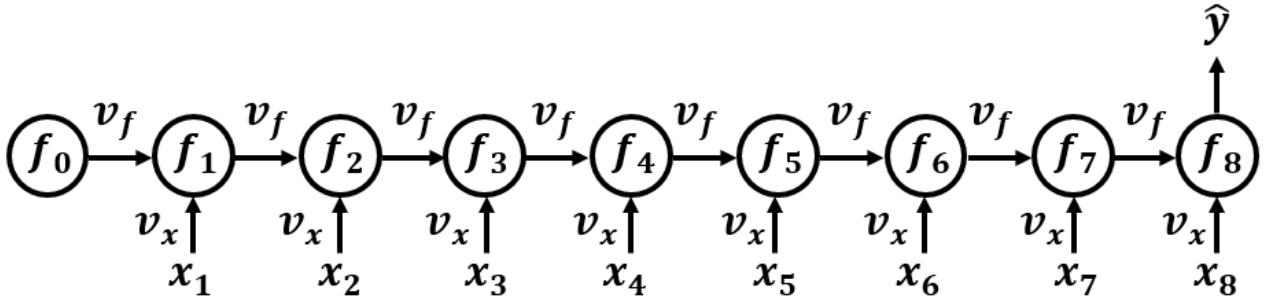


Figure 4: Sequence of states

where there is an abuse of notation from  $^{(t)}f$  to  $f_t$ .

Then, I can describe mathematically the above graph as follows

## 2.1 Forward propagation

$$\begin{aligned} f_t &= f_{t-1}v_f + x_tv_x \\ \hat{y} &= S_T \end{aligned} \tag{16}$$

## 2.2 Backpropagation

We define the error (loss) function as

$$E = \frac{1}{2} \sum_{i=1}^I (\hat{y}^{(i)} - y^{(i)})^2 \tag{17}$$

The parameters can be updated using the gradient descent method:

$$\begin{aligned} v_x &\leftarrow v_x - \alpha_x \frac{\partial E}{\partial v_x} \\ v_f &\leftarrow v_f - \alpha_f \frac{\partial E}{\partial v_f} \end{aligned} \quad (18)$$

Now, let us try to find  $\frac{\partial E}{\partial v_x}$  and  $\frac{\partial E}{\partial v_f}$ .

$$\frac{\partial E}{\partial v_x} = \sum_{t=1}^T \frac{\partial E}{\partial f_t} \frac{\partial f_t}{\partial v_x} \quad (19)$$

First, we compute  $\frac{\partial E}{\partial f_t}$  for  $t = 1 : 8$  as follows

$$\begin{aligned} \frac{\partial E}{\partial f_8} &= \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial f_8} = \sum_{i=1}^I (\hat{y}^{(i)} - y^{(i)}) \\ \frac{\partial E}{\partial f_7} &= \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial f_8} \frac{\partial f_8}{\partial f_7} = \sum_{i=1}^I (\hat{y}^{(i)} - y^{(i)}) v_f \\ &\vdots \\ \frac{\partial E}{\partial f_t} &= \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial f_8} \frac{\partial f_8}{\partial f_7} \dots \frac{\partial f_{t+1}}{\partial f_t} = \sum_{i=1}^I (\hat{y}^{(i)} - y^{(i)}) v_f^{(T-t)} \\ &\vdots \\ \frac{\partial E}{\partial f_1} &= \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial f_8} \frac{\partial f_8}{\partial f_7} \dots \frac{\partial f_2}{\partial f_1} = \sum_{i=1}^I (\hat{y}^{(i)} - y^{(i)}) v_f^{(T-1)} \end{aligned} \quad (20)$$

On the other hand,

$$\frac{\partial f_t}{\partial v_x} = x_t^{(i)} \quad (21)$$

Hence,

$$\begin{aligned} \frac{\partial E}{\partial v_x} &= \sum_{t=1}^T \sum_{i=1}^I (\hat{y}^{(i)} - y^{(i)}) v_f^{(T-t)} x_t^{(i)} \\ &= \sum_{t=1}^T \left\{ \left[ \sum_{i=1}^I (\hat{y}^{(i)} - y^{(i)}) x_t^{(i)} \right] v_f^{(T-t)} \right\} \end{aligned} \quad (22)$$

Then, we compute  $\frac{\partial E}{\partial v_f}$ .

$$\frac{\partial E}{\partial v_f} = \sum_{t=1}^T \frac{\partial E}{\partial f_t} \frac{\partial f_t}{\partial v_f} \quad (23)$$

From the previous steps, we already know that

$$\frac{\partial E}{\partial f_t} = \sum_{i=1}^I (\hat{y}^{(i)} - y^{(i)}) v_f^{(T-t)} \quad (24)$$

And,

$$\frac{\partial f_t}{\partial v_f} = f_{t-1}^{(i)}. \quad (25)$$

Hence,

$$\frac{\partial E}{\partial v_f} = \sum_{t=1}^T \left\{ \left[ \sum_{i=1}^I (\hat{y}^{(i)} - y^{(i)}) f_{t-1}^{(i)} \right] v_f^{(T-t)} \right\} \quad (26)$$

### 3 Resilient propagation

This is an algorithm proposed by Martin Riedmiller and Heinrich Brann.

#### 3.1 Motivation

From the example considered previously, we know that

$$\frac{\partial E}{\partial f_1} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial f_T} \frac{\partial f_T}{\partial f_{T-1}} \dots \frac{\partial f_t}{\partial f_{t-1}} \dots \frac{\partial f_2}{\partial f_1} = \sum_{i=1}^I (\hat{y}^{(i)} - y^{(i)}) v_f^{(T-1)} \quad (27)$$

Suppose now that  $T \gg 1$ , then as  $T \rightarrow \infty$ ,

$$\frac{\partial E}{\partial f_1} \rightarrow \begin{cases} 0, & \text{if } v_f < 1 \quad (\text{known as vanishing problem}) \\ \infty, & \text{if } v_f > 1 \quad (\text{known as explosion problem}) \end{cases} \quad (28)$$

Then, how can we solve this problem? A simple approach could be using the **resilient propagation**.

#### 3.2 Definition

The **resilient propagation** is similar to the backpropagation approach (in the sense of computing gradients), but it is different in that it does not use the actual values of the gradients. But, it is different in that it does not use the actual values of the gradients. But, it only uses the information coming from the signs of the gradients.

Hence, the new update rules are,

$$\begin{aligned} v_x &\leftarrow v_x - \text{sign} \left( \frac{\partial E}{\partial v_x} \right) \Delta_x \\ v_f &\leftarrow v_f - \text{sign} \left( \frac{\partial E}{\partial v_f} \right) \Delta_f \end{aligned} \quad (29)$$

where

$$\Delta_x \leftarrow \begin{cases} \Delta_x \cdot \eta_p, & \text{if } \text{sign} \left( \frac{\partial E}{\partial v_x} \right)_{\text{previous}} = \text{sign} \left( \frac{\partial E}{\partial v_x} \right)_{\text{current}} \\ \Delta_x \cdot \eta_n, & \text{if } \text{sign} \left( \frac{\partial E}{\partial v_x} \right)_{\text{previous}} \neq \text{sign} \left( \frac{\partial E}{\partial v_x} \right)_{\text{current}} \end{cases} \quad (30)$$

$$\Delta_f \leftarrow \begin{cases} \Delta_f \cdot \eta_p, & \text{if } \text{sign} \left( \frac{\partial E}{\partial v_f} \right)_{\text{previous}} = \text{sign} \left( \frac{\partial E}{\partial v_f} \right)_{\text{current}} \\ \Delta_f \cdot \eta_n, & \text{if } \text{sign} \left( \frac{\partial E}{\partial v_f} \right)_{\text{previous}} \neq \text{sign} \left( \frac{\partial E}{\partial v_f} \right)_{\text{current}} \end{cases} \quad (31)$$

where  $\eta_p = 1.2 > 1$ ,  $\eta_n = 0.5 < 1$ , and  $\Delta_{x_{\text{init}}} = \Delta_{f_{\text{init}}} = 0.001$ .

### 3.2.1 Why $\eta_n = 0.5 < 1$ ?

Every time the partial derivative of the corresponding weight changes its sign, this indicates that the last update was too big and the algorithm has jumped over a local minimum. Hence, the update value is decreased by  $\eta_n$  to reduce the update step size.

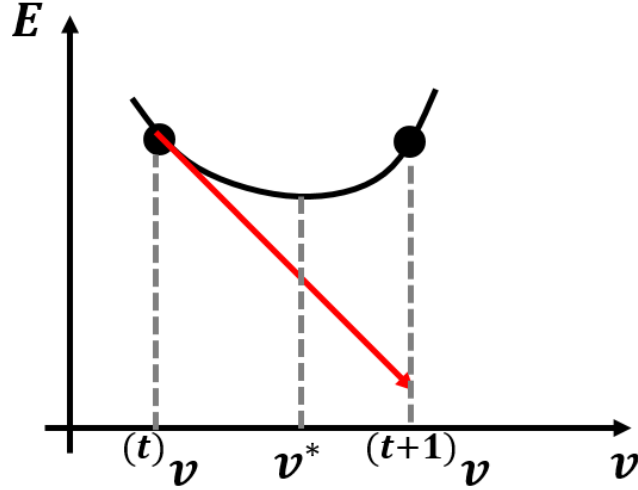


Figure 5: Resilient propagation with large step size

### 3.2.2 Why $\eta_p = 1.2 > 1$ ?

If the derivative retains its sign, the update value is slightly increased in order to accelerate the convergence in shallow regions.

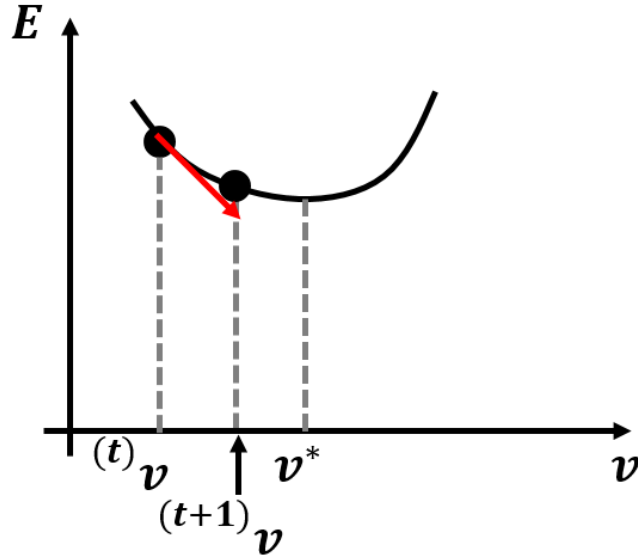


Figure 6: Resilient propagation with small step size



## 4 Gradient clipping

Another simple way to (possibly) get around the explosion problem is by clipping the gradients. If  $\|\frac{\partial E}{\partial V}\| > \eta$  for some  $\eta > 0$ , then clip the gradient as

$$\frac{\partial E}{\partial V} \longleftarrow \eta \frac{\frac{\partial E}{\partial V}}{\|\frac{\partial E}{\partial V}\|}.$$