

FASTRADE
FPGA-Accelerated Strategy & Trading for Enhanced Decisions

Weitao Lu, Wenbo Liu, Xiaolei Zhao, Yixuan Li

Content

- I. Introduction
- II. Hardware-Software Interfaces
- III. Hardware
- IV. Software
- V. Conclusion
- VI. Demo

Background

- **Raw Features:**
from Kaggle dataset

Date	# Open	# High	# Low	# Close	# Volume
2021-01-13 00:00:00-05:00	90.900001525878 9	103.0	90.010002136230 47	96.364997863769 53	25159000
2021-01-14 00:00:00-05:00	103.5	137.97999572753 906	101.0	114.94000244140 625	20858700

- **Factors:**

Price_Momentum	Volume_Factor	RSI	30d_Moving_Average
-0.034296688408046300	0.9885964411185060	43.91044444555740	111.4864995320640
-0.2044544936561990	1.1136913884395500	26.84116719152140	110.70316619873000
-0.2153845893012150	1.1066665475028100	17.19115889168140	109.86316630045600

- **Factor Model: Multi-Factor Model (5 Raw Features + 4 Factors)**

Reference:

- [Rendle, 2010] Steffen Rendle. Factorization machines. In ICDM, 2010.

Motivation

- **FPGAs are widely used at quant companies.**

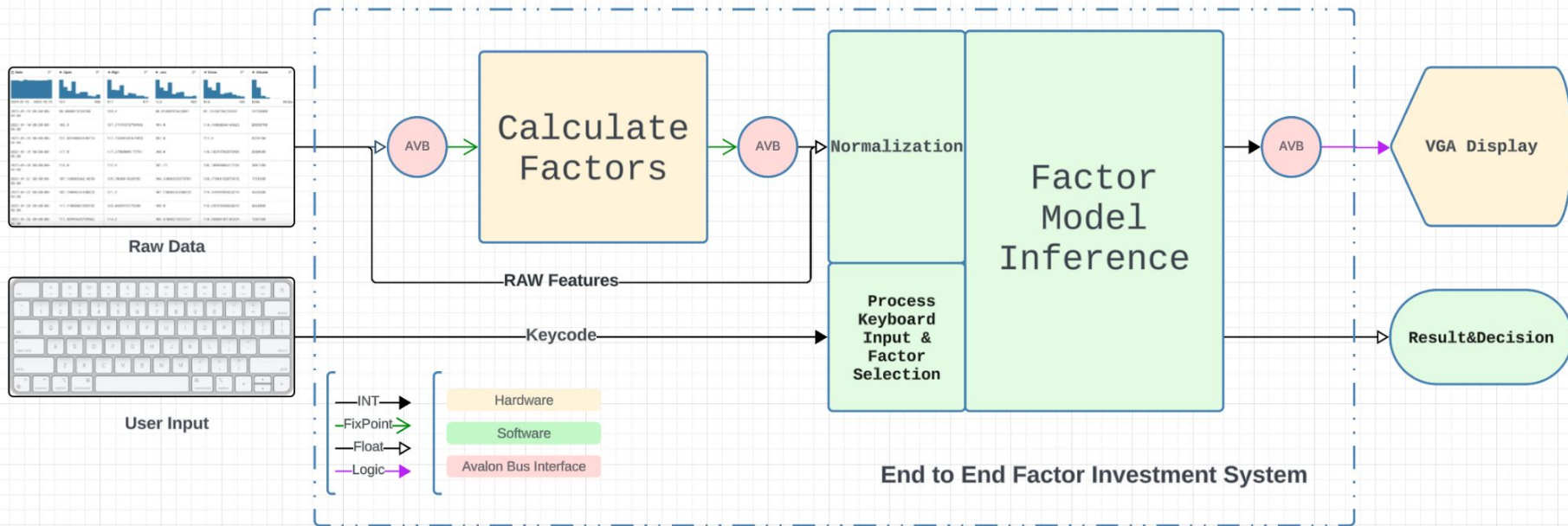
But quantitative researchers prefer Python.

We are interested to reproduce the investment pipeline on FPGA

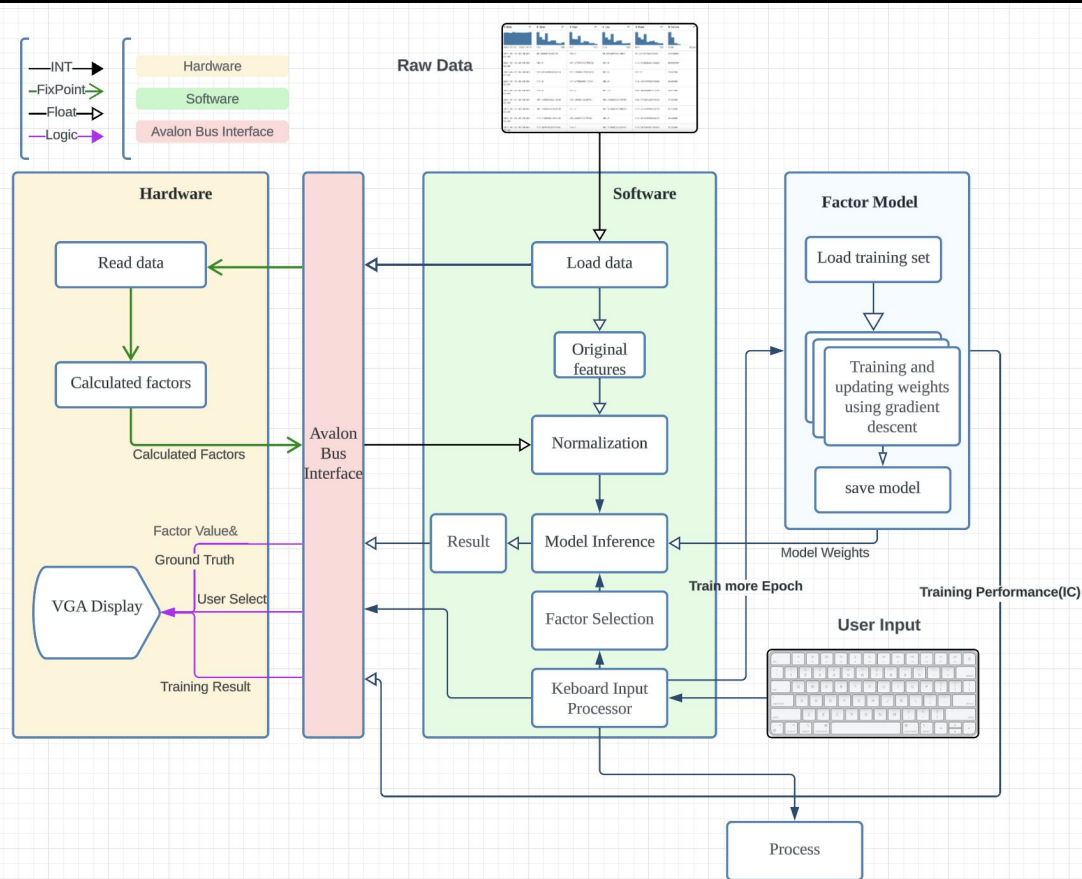
A Well Composed End-2-End Multi-Factor Investment Pipeline :

- **Components:Hardware Factor Calculation, Software model training/inference**
- **Well designed Hardware/Software data interfaces**
- **User-Friendly VGA display and Keyboard interaction.**
- **Verified hardware calculation waveform and acceleration improvement**

End2End System Overview



Detailed System Design



Hardware - Software Interfaces

The screenshot displays the Xilinx Component Editor interface for a project named 'soc_system.qsys'. The main window shows a table of components and their connections. The 'Connections' tab is active, showing a tree view of components and their interconnections. The 'Messages' window at the bottom shows two information messages regarding HPS Main PLL counter settings.

Use	Connections	Name	Description	Export	Clock	Base	End
<input checked="" type="checkbox"/>		clk_0	Clock Source	clk	exported		
<input checked="" type="checkbox"/>		clk_in	Clock Input	reset	clk_0		
<input checked="" type="checkbox"/>		clk_in_reset	Reset Input	Double-click to			
<input checked="" type="checkbox"/>		clk_reset	Reset Output	Double-click to			
<input checked="" type="checkbox"/>		h2f_user1_clock	Arria V/Cyclone V Hard Proc...	Double-click to			
<input checked="" type="checkbox"/>		hps_0	Conduit	hps_0 h2...			
<input checked="" type="checkbox"/>		hps_io	Conduit	hps			
<input checked="" type="checkbox"/>		h2f_reset	Reset Output	Double-click to	clk_0		
<input checked="" type="checkbox"/>		h2f_axi_clock	Clock Input	Double-click to	[h2f_axi_...		
<input checked="" type="checkbox"/>		h2f_axi_master	AXI Master	Double-click to	clk_0		
<input checked="" type="checkbox"/>		f2h_axi_slave	AXI Slave	Double-click to	[f2h_axi_...	#	
<input checked="" type="checkbox"/>		h2f_hw_axi_clock	Clock Input	Double-click to	clk_0		
<input checked="" type="checkbox"/>		h2f_hw_axi_master	AXI Master	Double-click to	[h2f_hw_...		
<input checked="" type="checkbox"/>		f2h_irq0	Interrupt Receiver	Double-click to		IRQ 0	IRQ 31
<input checked="" type="checkbox"/>		f2h_irq1	Interrupt Receiver	Double-click to		IRQ 0	IRQ 31
<input checked="" type="checkbox"/>		vga_ball_0	VGA Ball	Double-click to	clk_0		
<input checked="" type="checkbox"/>		clock	Clock Input	Double-click to	[clock]		
<input checked="" type="checkbox"/>		reset	Reset Input	Double-click to	[clock]	# 0x0000_0000	0x0000_003f
<input checked="" type="checkbox"/>		avalon_slave_0	Avalon Memory Mapped Slave	Double-click to	[clock]		
<input checked="" type="checkbox"/>		vga	Conduit	Double-click to			

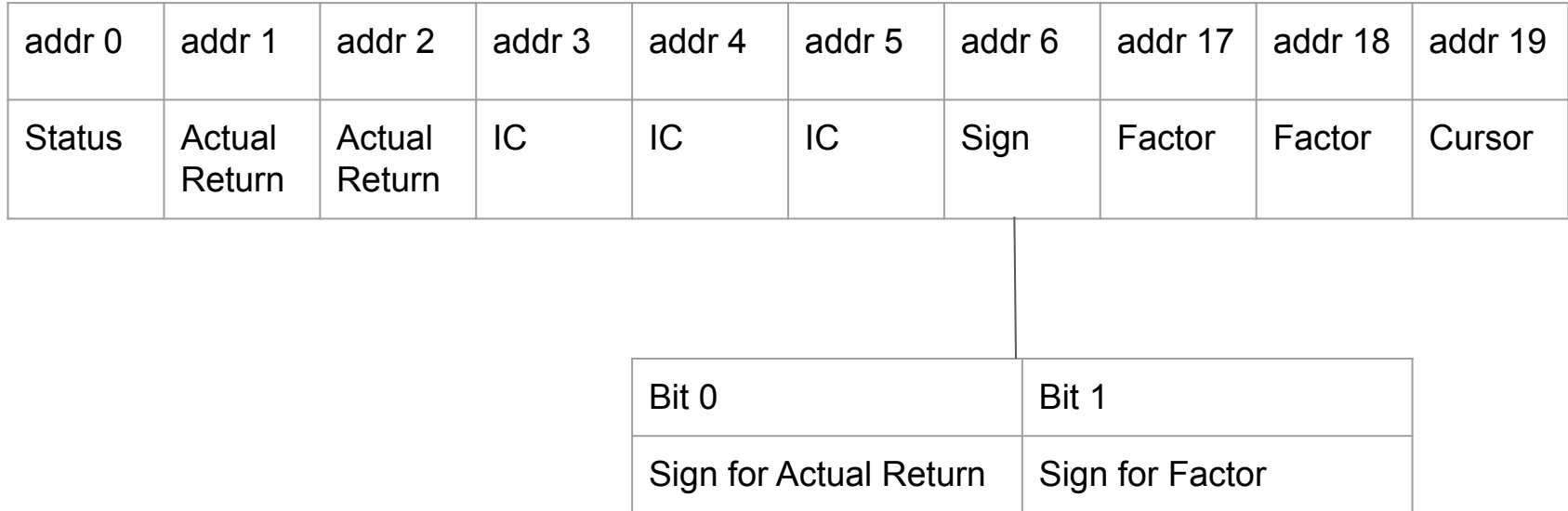
Messages

Type	Path	Message
Info		2 Info Messages
Info	soc_system.hps_0	HPS Main PLL counter settings: n = 0 m = 73
Info	soc_system.hps_0	HPS peripheral PLL counter settings: n = 0 m = 39

Component Editor - vga_ball_hw.tct*

```
File Templates Beta View
Component Type Block Symbol Files Parameters Signals & Interfaces
About Signals
Name
- avalon_slave_0 Avalon Memory Map
  - address [5] address
  - chipselect [1] chipselect
  - read [1] read
  - readdata [16] readdata
  - write [1] write
  - writedata [16] writedata
  <<add signal>>
- clock Clock Input
  - clk [1] clk
- reset Reset Input
  - reset [1] reset
  <<add signal>>
- vga Conduit
  - VGA_B [8] b
  - VGA_BLANK_n [1] blank_n
  - VGA_CLK [1] clk
  - VGA_G [8] g
  - VGA_HS [1] hs
  - VGA_R [8] r
  - VGA_SYNC_n [1] sync_n
  - VGA_VS [1] vs
  <<add signal>>
<<add interface>>
```


Address Map - VGA



Address Map - Calculator

Write to Hardware

addr 7	addr 8	addr 9	addr 10	addr 11
Open	High	Low	Close	Volume

Read from Hardware

addr 12	addr 13	addr 14	addr 15	addr 16
Momentum	Volume	RSI	MA	Done

Division Module

This module was designed for signed division of fixed-point numbers.

Inputs: Dividend, Divisor

Outputs: Quotient, Signal Completion, Busy, Divide by zero flag

Used Gaussian Rounding to implement the module and handled some special cases, including divided by zero.

Calculator Module

This module was designed to calculate the factors.

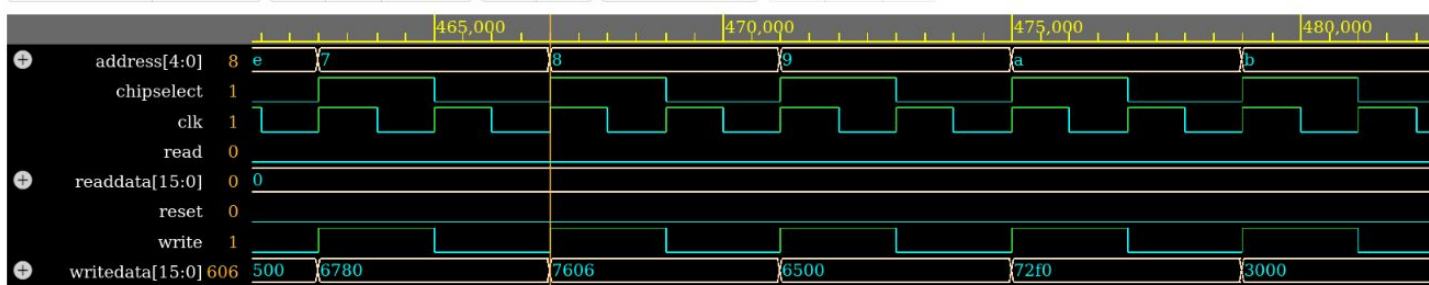
Use a 2D register array to store the original data and the calculation results(factors).

The hardware will read or write data according to the signals.

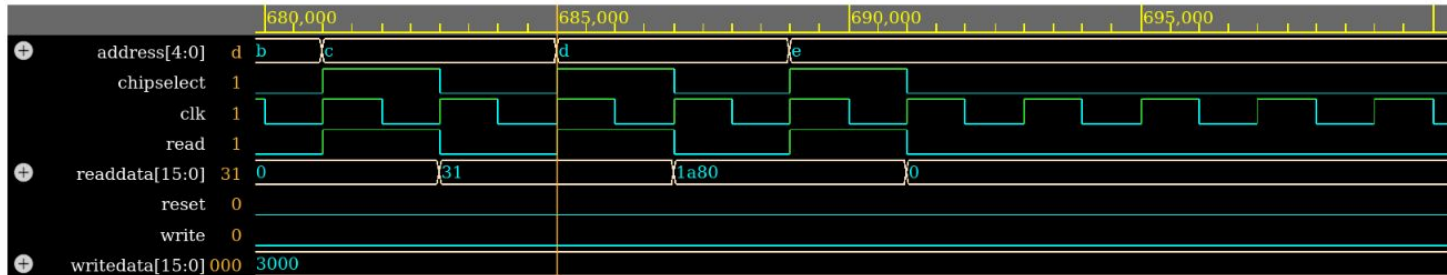
Only when all the five base data have been received, the module will start to calculate the factors.

Only when all the factors of one day have been calculated, the software part can read the results.

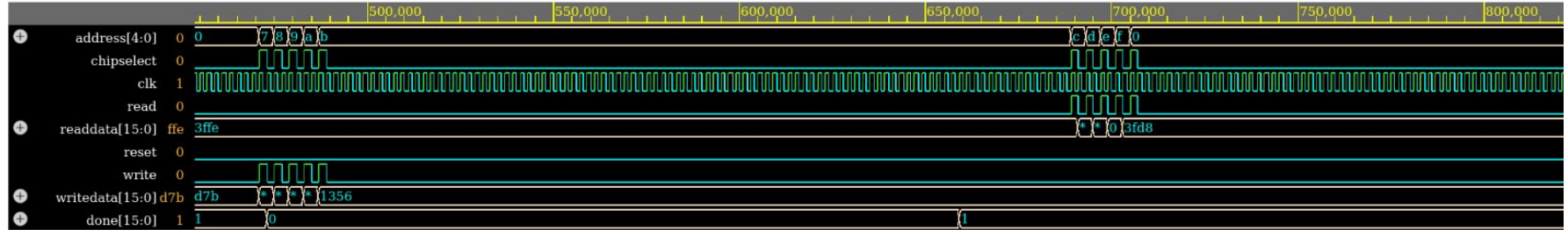
Calculator



Note: To revert to EDWave, click on a new browser window, not that option on your profile page.



Calculator

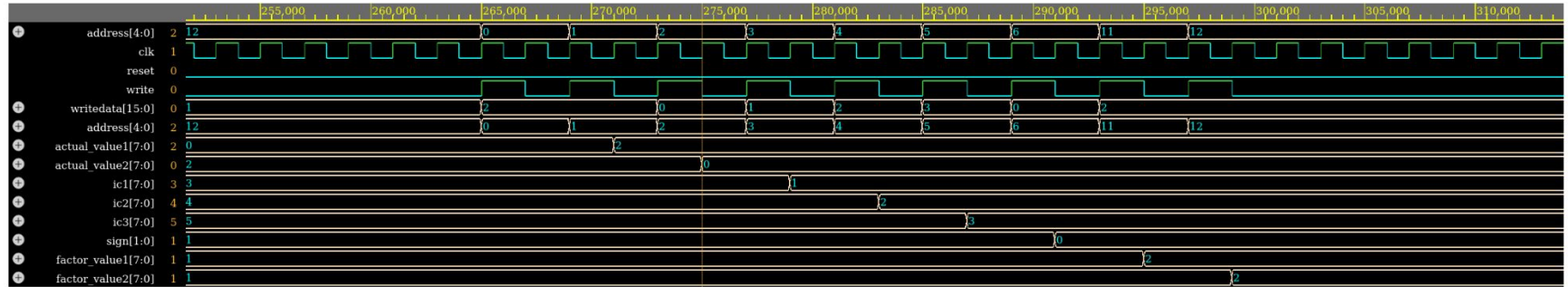


- Resolution: 680 * 480
 - 80 columns, 60 rows
- 50 character bitmaps: 8 * 8 bits/each
- Efficient storage(only store bitmaps)

```
char_a[0] = 8'b00011000; // **
char_a[1] = 8'b00100100; // * *
char_a[2] = 8'b01000010; // * *
char_a[3] = 8'b01111110; // *****
char_a[4] = 8'b01000010; // * *
char_a[5] = 8'b01000010; // * *
char_a[6] = 8'b01000010; // * *
char_a[7] = 8'b01000010; // * *

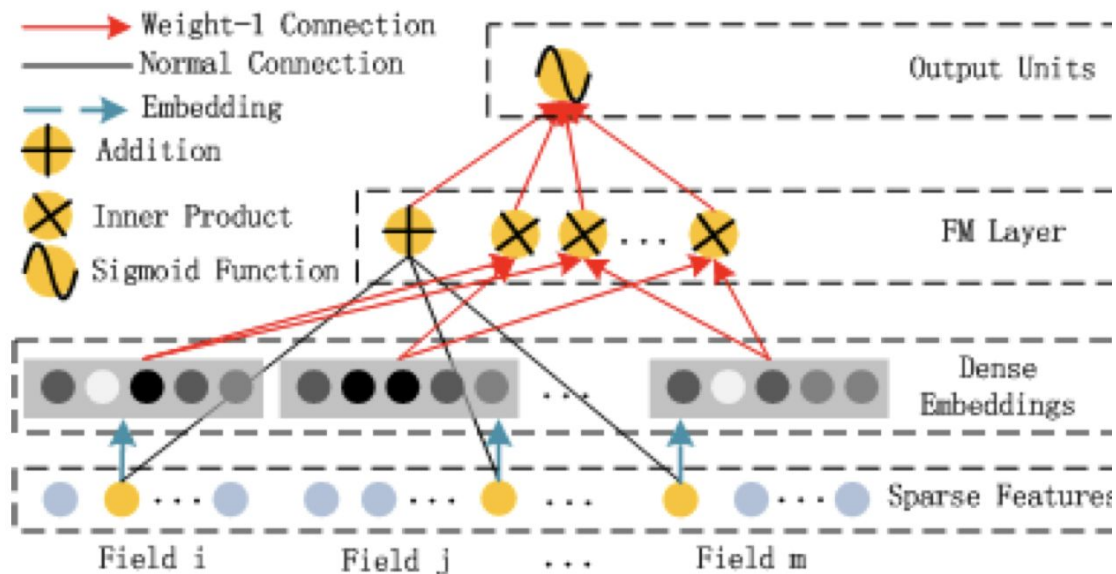
char_b[0] = 8'b01111100; // *****
char_b[1] = 8'b01000010; // * *
char_b[2] = 8'b01000010; // * *
char_b[3] = 8'b01111100; // *****
char_b[4] = 8'b01000010; // * *
char_b[5] = 8'b01000010; // * *
char_b[6] = 8'b01000010; // * *
char_b[7] = 8'b01111100; // *****

char_c[0] = 8'b00111100; // ****
char_c[1] = 8'b01000010; // * *
char_c[2] = 8'b01000000; // *
char_c[3] = 8'b01000000; // *
char_c[4] = 8'b01000000; // *
char_c[5] = 8'b01000000; // *
char_c[6] = 8'b01000010; // * *
char_c[7] = 8'b00111100; // ****
```



Note: To revert to EPWave online in a new browser window, click that online on your profile page.

- Factor Model



Software Problems and Solution

Everything works fine on Python, yet hard to program in C.

Lack of machine learning libs: -> Reproduced every from scratch by stand C libs.

```
// matrix.h
#ifndef MATRIX_H
#define MATRIX_H

typedef struct {
    double **data;
    int rows;
    int cols;
} Matrix;

Matrix create_matrix(int rows, int cols);
void free_matrix(Matrix m);
Matrix multiply_matrices(Matrix a, Matrix b);
void elementwise_multiply(Matrix a, Matrix b, Matrix result);
double sum_matrix_elements(Matrix m);
void init_random_normal(Matrix m, double mean, double stddev);
double sigmoid(double x);
void compute_means(Matrix features, double *means);
void compute_stddevs(Matrix features, double *means, double
*stddevs);
void standardize_features(Matrix features, double *means, double
*stddevs);
double pearson_correlation(double *x, double **y, int n);
#endif
```

```
// fm_model.h
#ifndef FM_MODEL_H
#define FM_MODEL_H

#include "matrix.h"
#include <stdbool.h>

typedef struct {
    double bias; // Scalar bias term w0 in the model
    Matrix weights; // Weight vector w in the model
    Matrix factors; // Factorization matrix V in the model
} FMModel;

void fit(FMModel *model, Matrix X, Matrix y, int feature_potential,
double alpha, int iter, int batch_size, double decay_rate);
double *predict(FMModel *model, Matrix X);
double *predict_active(FMModel *model, Matrix X, Matrix Ref ,bool
*active_features);
#endif
```

Software Problems and Solution

Everything works fine on Python, yet hard to program in C.

Overflow during gradient descent: -> Improved Training Strategy

- **Various Feature Scale, too large data:**
 - Zscore Normalization

- **Overflow caused by gradient accumulation:**
 - Mini-batch Gradient Descent
 - Learning Rate decay

Evaluation

- Hardware(Python runs on Google Colab)

Calculation Time using Python(Colab)	Calculation Time using Python(M1)	Calculation Time using SV
0.0050289 s	0.0014479 s	0.0018 s

- Software

Python Training /20epoch(Colab)	Python Inference Time(Colab)	Python Training /20epoch(M1)	Python Inference Time(M1)	On Board Training /20epoch	On Board Inference Time
1.103962 s	0.000150 s	0.225156 s	0.000026 s	0.130262 s	0.000035 s

*M1 refers to Apple M1 Pro and Colab refers to Google Colab's CPU

- **Blue** - current place of cursor
- **Green** - selected features
- Train - train the model
- Test - prediction

- Press ↑ and ↓ to move cursor
- Press enter to select/unselect feature



Thank You!

TRANSCENDING DISCIPLINES, TRANSFORMING LIVES