You are a senior chip design engineer who is an expert with Verilog HDL, designing AI accelerators for SoCs. I want to join efabless's 4 AI contest. The challenge aims to leverage generative AI to develop an open-source hardware accelerator designed explicitly for Keyword Spotting (KWS) applications on the Caravel System-on-Chip. I used the "Integer-Only Approximated MFCC for Ultra-Low Power Audio NN Processing on Multi-Core MCUs". To create the mfcc. The KWS uses the chose the CNN model that has been widely used for KWS, the CNN-KWS model, also known as the "Hello Edge" model. This model was introduced in the paper "Hello Edge: Keyword Spotting on Microcontrollers" by Zhang et al. You will help me debug the KWS machine learning and the audio features extractor and create an energy-efficient KWS accelerator that seamlessly integrates into the Caravel SoC environment. You will also help me add feature to the codebase. My top module is cnn_kws_accel "module cnn_kws_accel #( parameter NUM_KEYWORDS = 10, parameter MFCC_FEATURES = 40, parameter ACTIV_BITS = 8, parameter FC1_INPUT_SIZE = (MFCC_FEATURES/2)*(CONV2_NUM_FILTERS), parameter FC1_OUTPUT_SIZE = 64, parameter FC2_INPUT_SIZE = 64, parameter FC2_OUTPUT_SIZE = NUM_KEYWORDS, parameter CONV1_KERNEL_SIZE = 3, parameter CONV1_NUM_FILTERS = 8, parameter CONV2_KERNEL_SIZE = 3, parameter CONV2_NUM_FILTERS = 8 )( input wire clk, input wire rst_n, input wire [15:0] audio_in, input wire audio_valid, output reg [NUM_KEYWORDS-1:0] kws_result, output reg kws_valid, input wire [7:0] frame_size, input wire [7:0] frame_overlap, input wire [7:0] num_mfcc_coeffs, input wire [7:0] goertzel_coefs, output wire goertzel_coefs_start, output wire goertzel_coefs_valid, output wire goertzel_coefs_done, input wire [CONV1_NUM_FILTERS*CONV1_KERNEL_SIZE*CONV1_KERNEL_SIZE*ACTIV_BITS-1:0] conv1_weights, input wire [CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_biases, input wire conv1_load_weights, input wire conv1_load_biases, input wire [CONV2_NUM_FILTERS*CONV1_NUM_FILTERS*CONV2_KERNEL_SIZE*CONV2_KERNEL_SIZE*ACTIV_BITS-1:0] conv2_weights, input wire [CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_biases, input wire conv2_load_weights, input wire conv2_load_biases, input wire [FC1_OUTPUT_SIZE*FC1_INPUT_SIZE*ACTIV_BITS-1:0] fc1_weights, input wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_biases, input wire fc1_load_weights, input wire fc1_load_biases, input wire [FC2_OUTPUT_SIZE*FC2_INPUT_SIZE*ACTIV_BITS-1:0] fc2_weights, input wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_biases, input wire fc2_load_weights, input wire fc2_load_biases ); // MFCC module signals wire [15:0] mfcc_out; wire mfcc_valid; // CNN-KWS layers wire [MFCC_FEATURES*CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_out; wire conv1_valid; wire [MFCC_FEATURES*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_out; wire conv2_valid; wire [(MFCC_FEATURES/2)*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] maxpool_out; wire maxpool_valid; wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_out; wire fc1_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_out; wire fc2_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] softmax_out; wire softmax_valid; // MFCC module instantiation mfcc_accel mfcc ( .clk(clk), .rst(rst_n), // Adjusted to rst_n .audio_sample(audio_in), // Adjusted to audio_sample .sample_valid(audio_valid), // Adjusted to sample_valid .mfcc_feature(mfcc_out), // Adjusted to mfcc_feature .mfcc_valid(mfcc_valid) // Adjusted to mfcc_valid ); // Convolutional layer 1 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(1), .KERNEL_SIZE(CONV1_KERNEL_SIZE), .NUM_FILTERS(CONV1_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv1 ( .clk(clk), .rst_n(rst_n), .data_in(mfcc_out), .data_valid(mfcc_valid), .data_out(conv1_out), .data_out_valid(conv1_valid), .weights_in(conv1_weights),

.biases_in(conv1_biases), .load_weights(conv1_load_weights), .load_biases(conv1_load_biases) ); // Convolutional layer 2 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV1_NUM_FILTERS), .KERNEL_SIZE(CONV2_KERNEL_SIZE), .NUM_FILTERS(CONV2_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv2 ( .clk(clk), .rst_n(rst_n), .data_in(conv1_out), .data_valid(conv1_valid), .data_out(conv2_out), .data_out_valid(conv2_valid), .weights_in(conv2_weights), .biases_in(conv2_biases), .load_weights(conv2_load_weights), .load_biases(conv2_load_biases) ); maxpool2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV2_NUM_FILTERS), .KERNEL_SIZE(2), .STRIDE(2), .ACTIV_BITS(ACTIV_BITS) ) maxpool ( .clk(clk), .rst_n(rst_n), .data_in(conv2_out), .data_valid(conv2_valid), .data_out(maxpool_out), .data_out_valid(maxpool_valid) ); fully_connected #( .INPUT_SIZE(FC1_INPUT_SIZE), .OUTPUT_SIZE(FC1_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc1 ( .clk(clk), .rst_n(rst_n), .data_in(maxpool_out), .data_valid(maxpool_valid), .data_out(fc1_out), .data_out_valid(fc1_valid), .weights_in(fc1_weights[FC1_OUTPUT_SIZE*FC1_INPUT_SIZE*ACTIV_BITS-1:0]), .biases_in(fc1_biases), .load_weights(fc1_load_weights), .load_biases(fc1_load_biases) ); // Fully connected layer 2 (output layer) fully_connected #( .INPUT_SIZE(FC2_INPUT_SIZE), .OUTPUT_SIZE(FC2_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc2 ( .clk(clk), .rst_n(rst_n), .data_in(fc1_out), .data_valid(fc1_valid), .data_out(fc2_out), .data_out_valid(fc2_valid), .weights_in(fc2_weights), .biases_in(fc2_biases), .load_weights(fc2_load_weights), .load_biases(fc2_load_biases) ); // Softmax layer softmax #( .INPUT_SIZE(NUM_KEYWORDS), .ACTIV_BITS(ACTIV_BITS) ) softmax ( .clk(clk), .rst_n(rst_n), .data_in(fc2_out), .data_valid(fc2_valid), .data_out(softmax_out), .data_out_valid(softmax_valid) ); // Output assignment always @(posedge clk or negedge rst_n) begin if (!rst_n) begin kws_result <= 'b0; kws_valid <= 1'b0; end else begin kws_result <= softmax_out[NUM_KEYWORDS-1:0]; kws_valid <= softmax_valid; end end endmodule ". Submodules "module conv2d #( parameter INPUT_WIDTH = 40, parameter INPUT_HEIGHT = 1, parameter INPUT_CHANNELS = 1, parameter KERNEL_SIZE = 3, parameter NUM_FILTERS = 8, parameter PADDING = 1, parameter ACTIV_BITS = 8 ) ( input wire clk, input wire rst_n, input wire [INPUT_WIDTH*INPUT_HEIGHT*INPUT_CHANNELS*ACTIV_BITS-1:0] data_in, input wire data_valid, output reg [INPUT_WIDTH*INPUT_HEIGHT*NUM_FILTERS*ACTIV_BITS-1:0] data_out, output reg data_out_valid, input wire [NUM_FILTERS*INPUT_CHANNELS*KERNEL_SIZE*KERNEL_SIZE*ACTIV_BITS-1:0] weights_in, input wire [NUM_FILTERS*ACTIV_BITS-1:0] biases_in, input wire load_weights, input wire load_biases ); // Declare weights and biases reg [ACTIV_BITS-1:0] weights [0:NUM_FILTERS-1][0:INPUT_CHANNELS-1][0:KERNEL_SIZE-1][0:KERNEL_SIZE-1]; reg [ACTIV_BITS-1:0] biases [0:NUM_FILTERS-1]; // Declare internal signals reg [ACTIV_BITS-1:0] input_buffer [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:INPUT_CHANNELS-1]; reg [2*ACTIV_BITS-1:0] conv_result [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:NUM_FILTERS-1]; reg [ACTIV_BITS-1:0] relu_result [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:NUM_FILTERS-1]; // Load weights and biases integer i_load, j_load, k_load, l_load; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset weights and biases for (i_load = 0; i_load < NUM_FILTERS; i_load = i_load + 1) begin for (j_load = 0; j_load < INPUT_CHANNELS; j_load = j_load + 1) begin for (k_load = 0; k_load < KERNEL_SIZE; k_load = k_load + 1) begin for (l_load = 0; l_load < KERNEL_SIZE; l_load = l_load + 1) begin weights[i_load][j_load][k_load][l_load] <= 0; end end end biases[i_load] <= 0; end end else begin // Load weights when load_weights is asserted if (load_weights) begin for (i_load = 0; i_load < NUM_FILTERS; i_load = i_load + 1) begin for (j_load = 0; j_load < INPUT_CHANNELS; j_load = j_load + 1) begin for (k_load = 0; k_load < KERNEL_SIZE; k_load = k_load + 1) begin for (l_load = 0; l_load < KERNEL_SIZE; l_load = l_load + 1) begin weights[i_load][j_load][k_load][l_load] <= weights_in[(i_load*INPUT_CHANNELS*KERNEL_SIZE*KERNEL_SIZE +

j_load*KERNEL_SIZE*KERNEL_SIZE + k_load*KERNEL_SIZE + l_load)*ACTIV_BITS +: ACTIV_BITS]; end end end end end // Load biases when load_biases is asserted if (load_biases) begin for (i_load = 0; i_load < NUM_FILTERS; i_load = i_load + 1) begin biases[i_load] <= biases_in[i_load*ACTIV_BITS +: ACTIV_BITS]; end end end end // Convolution operation integer m_conv, n_conv, p_conv, q_conv, i_conv, j_conv, k_conv; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset internal signals and output for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < INPUT_CHANNELS; k_conv = k_conv + 1) begin input_buffer[i_conv][j_conv][k_conv] <= 0; end for (m_conv = 0; m_conv < NUM_FILTERS; m_conv = m_conv + 1) begin conv_result[i_conv][j_conv][m_conv] <= 0; relu_result[i_conv][j_conv][m_conv] <= 0; end end end data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Shift input data into buffer for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < INPUT_CHANNELS; k_conv = k_conv + 1) begin if (j_conv < INPUT_WIDTH - 1) begin input_buffer[i_conv][j_conv][k_conv] <= input_buffer[i_conv][j_conv+1][k_conv]; end else begin input_buffer[i_conv][j_conv][k_conv] <= data_in[i_conv*INPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS + j_conv*INPUT_CHANNELS*ACTIV_BITS + k_conv*ACTIV_BITS +: ACTIV_BITS]; end end end end // Perform convolution for (m_conv = 0; m_conv < INPUT_HEIGHT; m_conv = m_conv + 1) begin for (n_conv = 0; n_conv < INPUT_WIDTH; n_conv = n_conv + 1) begin for (p_conv = 0; p_conv < NUM_FILTERS; p_conv = p_conv + 1) begin conv_result[m_conv][n_conv][p_conv] = {{(2*ACTIV_BITS-ACTIV_BITS){1'b0}}, biases[p_conv]}; for (q_conv = 0; q_conv < INPUT_CHANNELS; q_conv = q_conv + 1) begin for (i_conv = 0; i_conv < KERNEL_SIZE; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < KERNEL_SIZE; j_conv = j_conv + 1) begin if (m_conv + i_conv - PADDING >= 0 && m_conv + i_conv - PADDING < INPUT_HEIGHT && n_conv + j_conv - PADDING >= 0 && n_conv + j_conv - PADDING < INPUT_WIDTH) begin conv_result[m_conv][n_conv][p_conv] = conv_result[m_conv][n_conv][p_conv] + weights[p_conv][q_conv][i_conv][j_conv] * input_buffer[m_conv + i_conv - PADDING][n_conv + j_conv - PADDING][q_conv]; end end end end end end end // Apply ReLU activation for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < NUM_FILTERS; k_conv = k_conv + 1) begin relu_result[i_conv][j_conv][k_conv] = (conv_result[i_conv][j_conv][k_conv][2*ACTIV_BITS-1] == 0) ? conv_result[i_conv][j_conv][k_conv][ACTIV_BITS-1:0] : 0; end end end // Assign output for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < NUM_FILTERS; k_conv = k_conv + 1) begin data_out[i_conv*INPUT_WIDTH*NUM_FILTERS*ACTIV_BITS + j_conv*NUM_FILTERS*ACTIV_BITS + k_conv*ACTIV_BITS +: ACTIV_BITS] <= relu_result[i_conv][j_conv][k_conv]; end end end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule", "module fully_connected #( parameter INPUT_SIZE = 160, parameter OUTPUT_SIZE = 64, parameter ACTIV_BITS = 8 ) ( input wire clk, input wire rst_n, input wire [INPUT_SIZE*ACTIV_BITS-1:0] data_in, input wire data_valid, output reg [OUTPUT_SIZE*ACTIV_BITS-1:0] data_out, output reg data_out_valid, input wire [OUTPUT_SIZE*INPUT_SIZE*ACTIV_BITS-1:0] weights_in, input wire [OUTPUT_SIZE*ACTIV_BITS-1:0] biases_in, input wire load_weights, input wire load_biases ); // Declare weights and biases reg [ACTIV_BITS-1:0] weights [0:OUTPUT_SIZE-1][0:INPUT_SIZE-1]; reg [ACTIV_BITS-1:0] biases [0:OUTPUT_SIZE-1]; // Declare internal signals reg [2*ACTIV_BITS-1:0] acc_result [0:OUTPUT_SIZE-1]; reg [ACTIV_BITS-1:0] relu_result [0:OUTPUT_SIZE-1]; // Load weights and biases integer i_load, j_load; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset weights and biases for (i_load = 0; i_load < OUTPUT_SIZE; i_load = i_load + 1) begin for (j_load = 0; j_load < INPUT_SIZE; j_load = j_load + 1) begin weights[i_load][j_load] <= 0; end biases[i_load]

<= 0; end end else begin // Load weights when load_weights is asserted if (load_weights) begin for (i_load = 0; i_load < OUTPUT_SIZE; i_load = i_load + 1) begin for (j_load = 0; j_load < INPUT_SIZE; j_load = j_load + 1) begin weights[i_load][j_load] <= weights_in[(i_load*INPUT_SIZE + j_load)*ACTIV_BITS +: ACTIV_BITS]; end end end // Load biases when load_biases is asserted if (load_biases) begin for (i_load = 0; i_load < OUTPUT_SIZE; i_load = i_load + 1) begin biases[i_load] <= biases_in[i_load*ACTIV_BITS +: ACTIV_BITS]; end end end end // Fully connected layer operation integer i_fc, j_fc; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset internal signals and output for (i_fc = 0; i_fc < OUTPUT_SIZE; i_fc = i_fc + 1) begin acc_result[i_fc] <= 0; relu_result[i_fc] <= 0; end data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Perform matrix multiplication for (i_fc = 0; i_fc < OUTPUT_SIZE; i_fc = i_fc + 1) begin acc_result[i_fc] = {{(2*ACTIV_BITS-ACTIV_BITS){1'b0}}, biases[i_fc]}; for (j_fc = 0; j_fc < INPUT_SIZE; j_fc = j_fc + 1) begin acc_result[i_fc] = acc_result[i_fc] + weights[i_fc][j_fc] * data_in[j_fc*ACTIV_BITS +: ACTIV_BITS]; end end // Apply ReLU activation for (i_fc = 0; i_fc < OUTPUT_SIZE; i_fc = i_fc + 1) begin relu_result[i_fc] <= (acc_result[i_fc][2*ACTIV_BITS-1] == 0) ? acc_result[i_fc][ACTIV_BITS-1:0] : 0; end // Assign output for (i_fc = 0; i_fc < OUTPUT_SIZE; i_fc = i_fc + 1) begin data_out[i_fc*ACTIV_BITS +: ACTIV_BITS] <= relu_result[i_fc]; end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule", " module maxpool2d #( parameter INPUT_WIDTH = 40, parameter INPUT_HEIGHT = 1, parameter INPUT_CHANNELS = 8, parameter KERNEL_SIZE = 2, parameter STRIDE = 2, parameter ACTIV_BITS = 8 ) ( input wire clk, input wire rst_n, input wire [INPUT_WIDTH*INPUT_HEIGHT*INPUT_CHANNELS*ACTIV_BITS-1:0] data_in, input wire data_valid, output reg [(INPUT_WIDTH/STRIDE)*INPUT_CHANNELS*ACTIV_BITS-1:0] data_out, output reg data_out_valid ); localparam OUTPUT_WIDTH = INPUT_WIDTH / STRIDE; localparam OUTPUT_HEIGHT = INPUT_HEIGHT / STRIDE; // Declare internal signals reg [ACTIV_BITS-1:0] input_buffer [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:INPUT_CHANNELS-1]; reg [ACTIV_BITS-1:0] max_value [0:INPUT_CHANNELS-1]; // Max pooling operation integer i, j, k, m, n; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset internal signals and output for (i = 0; i < INPUT_HEIGHT; i = i + 1) begin for (j = 0; j < INPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin input_buffer[i][j][k] <= 0; end end end data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Shift input data into buffer for (i = 0; i < INPUT_HEIGHT; i = i + 1) begin for (j = 0; j < INPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin if (j < INPUT_WIDTH - 1) begin input_buffer[i][j][k] <= input_buffer[i][j+1][k]; end else begin input_buffer[i][j][k] <= data_in[i*INPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS + j*INPUT_CHANNELS*ACTIV_BITS + k*ACTIV_BITS +: ACTIV_BITS]; end end end end // Perform max pooling for (i = 0; i < OUTPUT_HEIGHT; i = i + 1) begin for (j = 0; j < OUTPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin max_value[k] = input_buffer[i*STRIDE][j*STRIDE][k]; for (m = 0; m < KERNEL_SIZE; m = m + 1) begin for (n = 0; n < KERNEL_SIZE; n = n + 1) begin if (i*STRIDE + m < INPUT_HEIGHT && j*STRIDE + n < INPUT_WIDTH) begin max_value[k] = (input_buffer[i*STRIDE + m][j*STRIDE + n][k] > max_value[k]) ? input_buffer[i*STRIDE + m][j*STRIDE + n][k] : max_value[k]; end end end data_out[i*OUTPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS + j*INPUT_CHANNELS*ACTIV_BITS + k*ACTIV_BITS +: ACTIV_BITS] <= max_value[k]; end end end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule", "", "module softmax #( parameter INPUT_SIZE = 10, parameter ACTIV_BITS = 8 ) ( input wire clk, input wire rst_n, input wire [INPUT_SIZE*ACTIV_BITS-1:0] data_in, input wire data_valid, output reg [INPUT_SIZE*ACTIV_BITS-1:0] data_out, output reg data_out_valid ); // Declare internal signals reg [ACTIV_BITS-1:0] exp_values [0:INPUT_SIZE-1]; reg [2*ACTIV_BITS-1:0] sum_exp; reg [ACTIV_BITS-1:0] softmax_values [0:INPUT_SIZE-1]; // Softmax activation function integer i; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset internal signals and output for (i = 0; i < INPUT_SIZE;

i = i + 1) begin exp_values[i] <= 0; softmax_values[i] <= 0; end sum_exp <= 0; data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Compute exponential values for (i = 0; i < INPUT_SIZE; i = i + 1) begin exp_values[i] <= data_in[i*ACTIV_BITS +: ACTIV_BITS]; end // Compute sum of exponential values sum_exp <= 0; for (i = 0; i < INPUT_SIZE; i = i + 1) begin sum_exp <= sum_exp + {{ACTIV_BITS{1'b0}}, exp_values[i]}; end // Compute softmax values for (i = 0; i < INPUT_SIZE; i = i + 1) begin softmax_values[i] <= (exp_values[i] << ACTIV_BITS) / sum_exp[2*ACTIV_BITS-1:ACTIV_BITS]; end // Assign output for (i = 0; i < INPUT_SIZE; i = i + 1) begin data_out[i*ACTIV_BITS +: ACTIV_BITS] <= softmax_values[i]; end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule". Help me debug the design.

ChatGPT

Debugging the provided Verilog design involves verifying the functionality of each module and ensuring proper integration. Here's a step-by-step approach to debug the design:

1. **Verify MFCC Module**: Ensure that the MFCC module correctly processes the audio input and outputs valid MFCC features.

2. **Verify Convolutional Layers**: Check the convolution operation, weight loading, and ReLU activation in the conv2d modules.

3. **Verify Max Pooling Layer**: Ensure the maxpool2d module performs correct max pooling.

4. **Verify Fully Connected Layers**: Check matrix multiplication and ReLU activation in the fully_connected modules.

5. **Verify Softmax Layer**: Ensure the softmax module computes the softmax values correctly.

6. **Integrate and Test the Entire System**: Verify the integration of all modules within the **cnn_kws_accel** top module.

Here are some common issues and debugging steps for each module:

**1. MFCC Module**

Check if the MFCC module correctly processes the audio input:

verilog

Copy code

mfcc_accel mfcc ( .clk(clk), .rst(rst_n), .audio_sample(audio_in), .sample_valid(audio_valid), .mfcc_feature(mfcc_out), .mfcc_valid(mfcc_valid) );

- **Test the MFCC Output**: Verify that **mfcc_out** and **mfcc_valid** are producing the expected output for a given audio input.

- **Simulation**: Simulate the MFCC module with known audio samples and compare the output against expected MFCC coefficients.

**2. Convolutional Layers**

Check weight loading and convolution operation:

verilog

Copy code

conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(1), .KERNEL_SIZE(CONV1_KERNEL_SIZE), .NUM_FILTERS(CONV1_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv1 ( .clk(clk), .rst_n(rst_n), .data_in(mfcc_out), .data_valid(mfcc_valid), .data_out(conv1_out), .data_out_valid(conv1_valid), .weights_in(conv1_weights), .biases_in(conv1_biases), .load_weights(conv1_load_weights), .load_biases(conv1_load_biases );

- **Weight Loading**: Ensure that weights are loaded correctly into the **weights** register array.

- **Convolution Operation**: Verify the convolution results against expected outputs. Check boundary conditions (padding).

- **Simulation**: Simulate with known input data and weights to validate the output.

## 3. Max Pooling Layer

Ensure the max pooling operation is correct:

verilog

Copy code

maxpool2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV2_NUM_FILTERS), .KERNEL_SIZE(2), .STRIDE(2), .ACTIV_BITS(ACTIV_BITS) ) maxpool ( .clk(clk), .rst_n(rst_n), .data_in(conv2_out), .data_valid(conv2_valid), .data_out(maxpool_out), .data_out_valid(maxpool_valid) );

- **Pooling Operation**: Check the max pooling results against expected outputs.

- **Simulation**: Simulate with known input data to validate the output.

## 4. Fully Connected Layers

Check matrix multiplication and ReLU activation:

verilog

Copy code

fully_connected #( .INPUT_SIZE(FC1_INPUT_SIZE), .OUTPUT_SIZE(FC1_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc1 ( .clk(clk), .rst_n(rst_n), .data_in(maxpool_out), .data_valid(maxpool_valid), .data_out(fc1_out), .data_out_valid(fc1_valid), .weights_in(fc1_weights[FC1_OUTPUT_SIZE*FC1_INPUT_SIZE*ACTIV_BITS-1:0]), .biases_in(fc1_biases), .load_weights(fc1_load_weights), .load_biases(fc1_load_biases );

- **Weight Loading**: Ensure that weights are loaded correctly.

- **Matrix Multiplication**: Verify the multiplication and accumulation results.

- **ReLU Activation**: Check if ReLU is applied correctly.

- **Simulation**: Simulate with known input data and weights to validate the output.

## 5. Softmax Layer

Check the softmax calculation:

verilog

```
softmax #( .INPUT_SIZE(NUM_KEYWORDS), .ACTIV_BITS(ACTIV_BITS) ) softmax ( .clk(clk),
.rst_n(rst_n), .data_in(fc2_out), .data_valid(fc2_valid), .data_out(softmax_out),
.data_out_valid(softmax_valid) );
```

- **Exponent Calculation**: Verify the exponential values are calculated correctly.

- **Normalization**: Ensure the output is normalized to sum to 1.

- **Simulation**: Simulate with known input data to validate the output.

**Integration Testing**

- **Stimulus Generation**: Generate test stimuli that cover a wide range of scenarios.

- **Testbench**: Create a testbench to simulate the entire system, providing input and checking outputs.

- **Functional Verification**: Ensure the system meets the functional requirements and the KWS application works as expected.

**Example Testbench**

Here's a simplified example of a testbench for the top module:

verilog

Copy code

```
module tb_cnn_kws_accel; reg clk; reg rst_n; reg [15:0] audio_in; reg audio_valid; reg [7:0]
frame_size; reg [7:0] frame_overlap; reg [7:0] num_mfcc_coeffs; reg [7:0] goertzel_coefs; reg
[CONV1_NUM_FILTERS*CONV1_KERNEL_SIZE*CONV1_KERNEL_SIZE*ACTIV_BITS-1:0]
conv1_weights; reg [CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_biases; reg conv1_load_weights;
reg conv1_load_biases; reg
[CONV2_NUM_FILTERS*CONV1_NUM_FILTERS*CONV2_KERNEL_SIZE*CONV2_KERNEL_SIZE*ACTIV_
BITS-1:0] conv2_weights; reg [CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_biases; reg
conv2_load_weights; reg conv2_load_biases; reg [FC1_OUTPUT_SIZE*FC1_INPUT_SIZE*ACTIV_BITS-
1:0] fc1_weights; reg [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_biases; reg fc1_load_weights; reg
fc1_load_biases; reg [FC2_OUTPUT_SIZE*FC2_INPUT_SIZE*ACTIV_BITS-1:0] fc2_weights; reg
[FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_biases; reg fc2_load_weights; reg fc2_load_biases; wire
[NUM_KEYWORDS-1:0] kws_result; wire kws_valid; // Instantiate the DUT (Device Under Test)
cnn_kws_accel dut ( .clk(clk), .rst_n(rst_n), .audio_in(audio_in), .audio_valid(audio_valid),
.kws_result(kws_result), .kws_valid(kws_valid), .frame_size(frame_size),
.frame_overlap(frame_overlap), .num_mfcc_coeffs(num_mfcc_coeffs),
.goertzel_coefs(goertzel_coefs), .conv1_weights(conv1_weights), .conv1_biases(conv1_biases),
.conv1_load_weights(conv1_load_weights), .conv1_load_biases(conv1_load_biases),
.conv2_weights(conv2_weights), .conv2_biases(conv2_biases),
.conv2_load_weights(conv2_load_weights), .conv2_load_biases(conv2_load_biases),
.fc1_weights(fc1_weights), .fc1_biases(fc1_biases), .fc1_load_weights(fc1_load_weights),
.fc1_load_biases(fc1_load_biases), .fc2_weights(fc2_weights), .fc2_biases(fc2_biases),
.fc2_load_weights(fc2_load_weights), .fc2_load_biases(fc2_load_biases) ); // Clock generation
always #5 clk = ~clk; initial begin // Initialize signals clk = 0; rst_n = 0; audio_in = 0; audio_valid = 0;
frame_size = 0; frame_overlap = 0; num_mfcc_coeffs = 0; goertzel_coefs = 0; conv1_weights = 0;
```

conv1_biases = 0; conv1_load_weights = 0; conv1_load_biases = 0; conv2_weights = 0; conv2_biases = 0; conv2_load_weights = 0; conv2_load_biases = 0; fc1_weights = 0; fc1_biases = 0; fc1_load_weights = 0; fc1_load_biases = 0; fc2_weights = 0; fc2_biases = 0; fc

I receive the following compiler errors "Resizing cell port cnn_kws_accel.conv1.data_in from 16 bits to 320 bits.".

ChatGPT

The error "Resizing cell port cnn_kws_accel.conv1.data_in from 16 bits to 320 bits" indicates a mismatch between the width of the data being provided to the **data_in** port of the **conv1** module and the width expected by that port. Specifically, the **mfcc_out** signal is 16 bits wide, while the **data_in** port of **conv1** expects a 320-bit wide signal (calculated based on the parameters **MFCC_FEATURES**, **INPUT_HEIGHT**, **INPUT_CHANNELS**, and **ACTIV_BITS**).

Let's walk through the necessary steps to correct this issue:

1. **Ensure Correct Data Widths**: Verify that each module's data ports match the expected widths.

2. **Adjust Data Widths or Interfaces**: If necessary, adjust the widths of the data signals or the module interfaces to ensure they match.

**Step-by-Step Solution**

1. **Verify the Widths of mfcc_out and conv1.data_in**

The **conv1** module's **data_in** port expects a width of **MFCC_FEATURES * INPUT_HEIGHT * INPUT_CHANNELS * ACTIV_BITS**. Given the default parameters:

- **MFCC_FEATURES** is 40.

- **INPUT_HEIGHT** is 1.

- **INPUT_CHANNELS** is 1.

- **ACTIV_BITS** is 8.

Therefore, **conv1.data_in** should be **40 * 1 * 1 * 8 = 320** bits wide.

The **mfcc_out** signal is only 16 bits wide, which is causing the mismatch.

2. **Adjust mfcc_out Signal Width**

Ensure that the **mfcc_out** signal is correctly sized to match the expected width for the **conv1** module. The **mfcc** module needs to output the correct number of features, each with the appropriate bit width.

3. **Check mfcc_accel Module**

Update the **mfcc_accel** module to output the correct width. Ensure it produces **MFCC_FEATURES** number of features, each **ACTIV_BITS** wide.

**Example Update to mfcc_accel Module**

verilog

Copy code

module mfcc_accel ( input wire clk, input wire rst_n, input wire [15:0] audio_sample, input wire sample_valid, output reg [MFCC_FEATURES*ACTIV_BITS-1:0] mfcc_feature, // Adjusted width output reg mfcc_valid ); // ... (implementation of MFCC feature extraction) endmodule

**Updated Top Module (cnn_kws_accel)**

Ensure that **mfcc_out** is now correctly sized and connects to **conv1** appropriately.

verilog

Copy code

module cnn_kws_accel #( parameter NUM_KEYWORDS = 10, parameter MFCC_FEATURES = 40, parameter ACTIV_BITS = 8, parameter FC1_INPUT_SIZE = (MFCC_FEATURES/2)*(CONV2_NUM_FILTERS), parameter FC1_OUTPUT_SIZE = 64, parameter FC2_INPUT_SIZE = 64, parameter FC2_OUTPUT_SIZE = NUM_KEYWORDS, parameter CONV1_KERNEL_SIZE = 3, parameter CONV1_NUM_FILTERS = 8, parameter CONV2_KERNEL_SIZE = 3, parameter CONV2_NUM_FILTERS = 8 )( input wire clk, input wire rst_n, input wire [15:0] audio_in, input wire audio_valid, output reg [NUM_KEYWORDS-1:0] kws_result, output reg kws_valid, input wire [7:0] frame_size, input wire [7:0] frame_overlap, input wire [7:0] num_mfcc_coeffs, input wire [7:0] goertzel_coefs, output wire goertzel_coefs_start, output wire goertzel_coefs_valid, output wire goertzel_coefs_done, input wire [CONV1_NUM_FILTERS*CONV1_KERNEL_SIZE*CONV1_KERNEL_SIZE*ACTIV_BITS-1:0] conv1_weights, input wire [CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_biases, input wire conv1_load_weights, input wire conv1_load_biases, input wire [CONV2_NUM_FILTERS*CONV1_NUM_FILTERS*CONV2_KERNEL_SIZE*CONV2_KERNEL_SIZE*ACTIV_BITS-1:0] conv2_weights, input wire [CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_biases, input wire conv2_load_weights, input wire conv2_load_biases, input wire [FC1_OUTPUT_SIZE*FC1_INPUT_SIZE*ACTIV_BITS-1:0] fc1_weights, input wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_biases, input wire fc1_load_weights, input wire fc1_load_biases, input wire [FC2_OUTPUT_SIZE*FC2_INPUT_SIZE*ACTIV_BITS-1:0] fc2_weights, input wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_biases, input wire fc2_load_weights, input wire fc2_load_biases ); // MFCC module signals wire [MFCC_FEATURES*ACTIV_BITS-1:0] mfcc_out; // Adjusted width wire mfcc_valid; // CNN-KWS layers wire [MFCC_FEATURES*CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_out; wire conv1_valid; wire [MFCC_FEATURES*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_out; wire conv2_valid; wire [(MFCC_FEATURES/2)*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] maxpool_out; wire maxpool_valid; wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_out; wire fc1_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_out; wire fc2_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] softmax_out; wire softmax_valid; // MFCC module instantiation mfcc_accel mfcc ( .clk(clk), .rst(rst_n), .audio_sample(audio_in), .sample_valid(audio_valid), .mfcc_feature(mfcc_out), .mfcc_valid(mfcc_valid) ); // Convolutional layer 1 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(1), .KERNEL_SIZE(CONV1_KERNEL_SIZE), .NUM_FILTERS(CONV1_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv1 ( .clk(clk), .rst_n(rst_n), .data_in(mfcc_out), // Corrected width .data_valid(mfcc_valid), .data_out(conv1_out), .data_out_valid(conv1_valid), .weights_in(conv1_weights), .biases_in(conv1_biases), .load_weights(conv1_load_weights), .load_biases(conv1_load_biases) ); // (The rest of the module instantiations and always block remain unchanged) endmodule

**Test and Simulate**

- **Simulate the MFCC module**: Verify that it outputs the correct width for **mfcc_out**.

- **Integrate and simulate the top module**: Ensure that the data flows correctly from the MFCC module through the convolutional layers, max pooling, fully connected layers, and softmax layer.

By ensuring that the widths match across the different modules, the error "Resizing cell port cnn_kws_accel.conv1.data_in from 16 bits to 320 bits" should be resolved.

Problem not fixed. Recalculate port and signal widths according to requirements. MFCC outputs 16bit Output: MFCC feat[M] as INT16 Q4 (16-bit integers with 4 fractional bits). Paper abstract "1 Introduction Deep learning algorithms have evolved to a stage where they have surpassed human accuracies in a variety of cognitive tasks including image classification [1] and conversational speech recognition [2]. Motivated by the recent breakthroughs in deep learning based speech recognition technologies, speech is increasingly becoming a more natural way to interact with consumer electronic devices, for example, Amazon Echo, Google Home and smart phones. However, always-on speech recognition is not energy-efficient and may also cause network congestion to transmit continuous audio stream from billions of these devices to the cloud. Furthermore, such a cloud based solution adds latency to the application, which hurts user experience. There are also privacy concerns when audio is continuously transmitted to the cloud. To mitigate these concerns, the devices first detect predefined keyword(s) such as "Alexa", "Ok Google", "Hey Siri", etc., which is commonly known as keyword spotting (KWS). Detection of keyword wakes up the device and then activates the full scale speech recognition either on device [3] or in the cloud. In some applications, the sequence of keywords can be used as voice commands to a smart device such as a voice-enabled light bulb. Since KWS system is always-on, it should have very low power consumption to maximize battery life. On the other hand, the KWS system should detect the keywords with high accuracy and low latency, for best user experience. These conflicting system requirements make KWS an active area of research ever since its inception over 50 years ago [4]. Recently, with the renaissance of artificial neural networks in the form of deep learning algorithms, neural network (NN) based KWS has become very popular [5, 6, 7, 8]. 2Work was done while the author was an intern at Arm. arXiv:1711.07128v3 [cs.SD] 14 Feb 2018 Low power consumption requirement for keyword spotting systems make microcontrollers an obvious choice for deploying KWS in an always-on system. Microcontrollers are low-cost energy-efficient processors that are ubiquitous in our everyday life with their presence in a variety of devices ranging from home appliances, automobiles and consumer electronics to wearables. However, deployment of neural network based KWS on microcontrollers comes with following challenges: Limited memory footprint: Typical microcontroller systems have only tens to few hundred KB of memory available. The entire neural network model, including input/output, weights and activations, has to fit within this small memory budget. Limited compute resources: Since KWS is always-on, the real-time requirement limits the total number of operations per neural network inference. These microcontroller resource constraints in conjunction with the high accuracy and low latency requirements of KWS call for a resource-constrained neural network architecture exploration to find lean neural network structures suitable for KWS, which is the primary focus of our work. The main contributions of this work are as follows: • We first train the popular KWS neural net models from the literature [5, 6, 7, 8] on Google speech commands dataset [9] and compare them in terms of accuracy, memory footprint and number of operations per inference. • In addition, we implement a new KWS model using depth-wise separable convolutions and point-wise convolutions, inspired by the success of resource-efficient MobileNet [10] in computer vision. This model outperforms the other prior models in all aspects of accuracy, model size and number of operations. • Finally, we perform resource-constrained neural network architecture exploration and

present comprehensive comparison of different network architectures within a set of compute and memory constraints of typical microcontrollers. The code, model definitions and pretrained models are available at https://github.com/ARM-software/ML-KWS-for-MCU. 2 Background 2.1 Keyword Spotting (KWS) System A typical KWS system consists of a feature extractor and a neural network based classifier as shown in Fig. 1. First, the input speech signal of length L is framed into overlapping frames of length l with a stride s, giving a total of T = L–l s + 1 frames. From each frame, F speech features are extracted, generating a total of T × F features for the entire input speech signal of length L. Log-mel filter bank energies (LFBE) and Mel-frequency cepstral coefficients (MFCC) are the commonly used human-engineered speech features in deep learning based speech-recognition, that are adapted from traditional speech processing techniques. Feature extraction using LFBE or MFCC involves translating the time-domain speech signal into a set of frequency-domain spectral coefficients, which enables dimensionality compression of the input signal. The extracted speech feature matrix is fed into a classifier module, which generates the probabilities for the output classes. In a real-world scenario where keywords need to be identified from a continuous audio stream, a posterior handling module averages the output probabilities of each output class over a period of time, improving the overall confidence of the prediction. Figure 1: Keyword spotting pipeline. 2 Traditional speech recognition technologies for KWS use Hidden Markov Models (HMMs) and Viterbi decoding [11, 12]. While these approaches achieve reasonable accuracies, they are hard to train and are computationally expensive during inference. Other techniques explored for KWS include discriminative models adopting a large-margin problem formulation [13] or recurrent neural networks (RNN) [14]. Although these methods significantly outperform HMM based KWS in terms of accuracy, they suffer from large detection latency. KWS models using deep neural networks (DNN) based on fully-connected layers with rectified linear unit (ReLU) activation functions are introduced in [5], which outperforms the HMM models with a very small detection latency. Furthermore, low-rank approximation techniques are used to compress the DNN model weights achieving similar accuracy with less hardware resources [15, 16]. The main drawback of DNNs is that they ignore the local temporal and spectral correlation in the input speech features. In order to exploit these correlations, different variants of convolutional neural network (CNN) based KWS are explored in [6], which demonstrate higher accuracy than DNNs. The drawback of CNNs in modeling time varying signals (e.g. speech) is that they ignore long term temporal dependencies. Combining the strengths of CNNs and RNNs, convolutional recurrent neural network based KWS is investigated in [7] and demonstrate the robustness of the model to noise. While all the prior KWS neural networks are trained with cross entropy loss function, a max-pooling based loss function for training KWS model with long short-term memory (LSTM) is proposed in [8], which achieves better accuracy than the DNNs and LSTMs trained with cross entropy loss. Although many neural network models for KWS are presented in literature, it is difficult to make a fair comparison between them as they are all trained and evaluated on different proprietary datasets (e.g. "TalkType" dataset in [7], "Alexa" dataset in [8], etc.) with different input speech features and audio duration. Also, the primary focus of prior research has been to maximize the accuracy with a small memory footprint model, without explicit constraints of underlying hardware, such as limits on number of operations per inference. In contrast, this work is more hardware-centric and targeted towards neural network architectures that maximize accuracy on microcontroller devices. The constraints on memory and compute significantly limit the neural network parameters and the number of operations. 2.2 Microcontroller Systems A typical microcontroller system consists of a processor core, an on-chip SRAM block and an on-chip embedded flash. Table 1 shows some commercially available microcontroller development boards with Arm Cortex-M processor cores with different compute capabilities running at different frequencies (16 MHz to 216 MHz), consisting of a wide range of on-chip memory (SRAM: 8 KB to 320 KB; Flash: 128 KB to 1 MB). The program binary, usually preloaded into the non-volatile flash, is

loaded into the SRAM at startup and the processor runs the program with the SRAM as the main data memory. Therefore, the size of the SRAM limits the size of memory that the software can use. Other than the memory footprint, performance (i.e., operations per second) is also a constraining factor for running neural networks on microcontrollers. Most microcontrollers are designed for embedded applications with low cost and high energy-efficiency as the primary targets, and do not have high throughput for compute-intensive workloads such as neural networks. Some microcontrollers have integrated DSP instructions that can be useful for running neural network workloads. For example, Cortex-M4 and Cortex-M7 have integrated SIMD and MAC instructions that can be used to accelerate low-precision computation in neural networks. Arm MbedTM platform Processor Frequency SRAM Flash Mbed LPC11U24 Cortex-M0 48 MHz 8 KB 32 KB Nordic nRF51-DK Cortex-M0 16 MHz 32 KB 256 KB Mbed LPC1768 Cortex-M3 96 MHz 32 KB 512 KB Nucleo F103RB Cortex-M3 72 MHz 20 KB 128 KB Nucleo L476RG Cortex-M4 80 MHz 128 KB 1 MB Nucleo F411RE Cortex-M4 100 MHz 128 KB 512 KB FRDM-K64F Cortex-M4 120 MHz 256 KB 1 MB Nucleo F746ZG Cortex M7 216 MHz 320 KB 1 MB Table 1: Typical off the shelf Arm Cortex-M based microcontroller development platforms. 3 3 Neural Network Architectures for KWS This section gives an overview of all the different neural network architectures explored in this work including the deep neural network (DNN), convolutional neural network (CNN), recurrent neural network (RNN), convolutional recurrent neural network (CRNN) and depthwise separable convolutional neural network (DS-CNN). 3.1 Deep Neural Network (DNN) The DNN is a standard feed-forward neural network made of a stack of fully-connected layers and non-linear activation layers. The input to the DNN is the flattened feature matrix, which feeds into a stack of d hidden fully-connected layers each with n neurons. Typically, each fully-connected layer is followed by a rectified linear unit (ReLU) based activation function. At the output is a linear layer followed by a softmax layer generating the output probabilities of the k keywords, which are used for further posterior handling. 3.2 Convolutional Neural Network (CNN) One main drawback of DNN based KWS is that they fail to efficiently model the local temporal and spectral correlation in the speech features. CNNs exploit this correlation by treating the input time-domain and spectral-domain features as an image and performing 2-D convolution operations over it. The convolution layers are typically followed by batch normalization [17], ReLU based activation functions and optional max/average pooling layers, which reduce the dimensionality of the features. During inference, the parameters of batch normalization can be folded into the weights of the convolution layers. In some cases, a linear low-rank layer, which is simply a fully-connected layer without non-linear activation, is added in between the convolution layers and dense layers for the purpose of reducing parameters and accelerating training [18, 19]. 3.3 Recurrent Neural Network (RNN) RNNs have shown superior performance in many sequence modeling tasks, especially speech recognition [20], language modeling [21], translation [22], etc. RNNs not only exploit the temporal relation between the input signal, but also capture the long-term dependencies using "gating" mechanism. Unlike CNNs where input features are treated as 2-D image, RNNs operate for T time steps, where at each time step t the corresponding spectral feature vector $f_t \in R^F$ concatenated with the previous time step output $h_{t-1}$ is used as input to the RNN. Figure 2 shows the model architecture of a typical RNN model, where the RNN cell could be an LSTM cell [23, 24] or a gated recurrent unit (GRU) cell [25, 26]. Since the weights are reused across all the T time steps, the RNN models tend to have less number of parameters compared to the CNNs. Similar to batch normalization in CNNs, research show that applying layer normalization can be beneficial for training RNNs [27], in which the hidden states are normalized during each time step. 3.4 Convolutional Recurrent Neural Network (CRNN) Convolution recurrent neural network [7] is a hybrid of CNN and RNN, which takes advantages of both. It exploits the local temporal/spatial correlation using convolution layers and global temporal dependencies in the speech features using recurrent layers. As shown in Fig. 3, a CRNN model Figure 2: Model architecture of RNN. 4 Figure 3: Model

Architecture of CRNN. starts with a convolution layer, followed by an RNN to encode the signal and a dense fully-connected layer to map the information. Here, the recurrent layer is bi-directional [28] and has multiple stages, increasing the network learning capability. Gated recurrent units (GRU) [25] is used as the base cell for recurrent layers, as it uses fewer parameters than LSTMs and gave better convergence in our experiments. 3.5 Depthwise Separable Convolutional Neural Network (DS-CNN) Recently, depthwise separable convolution has been proposed as an efficient alternative to the standard 3-D convolution operation [29] and has been used to achieve compact network architectures in the area of computer vision [10, 30]. DS-CNN first convolves each channel in the input feature map with a separate 2-D filter and then uses pointwise convolutions (i.e. 1x1) to combine the outputs in the depth dimension. By decomposing the standard 3-D convolutions into 2-D convolutions followed by 1-D convolutions, depthwise separable convolutions are more efficient both in number of parameters and operations, which makes deeper and wider architecture possible even in the resource-constrained microcontroller devices. In this work, we adopt a depthwise separable CNN based on the implementation of MobileNet [10] as shown in Fig. 4. An average pooling followed by a fully-connected layer is used at the end to provide global interaction and reduce the total number of parameters in the final layer. 4 Experiments and Results We use the Google speech commands dataset [9] for the neural network architecture exploration experiments. The dataset consists of 65K 1-second long audio clips of 30 keywords, by thousands of different people, with each clip consisting of only one keyword. The neural network models are trained to classify the incoming audio into one of the 10 keywords - "Yes", "No", "Up", "Down", "Left", "Right", "On", "Off", "Stop", "Go", along with "silence" (i.e. no word spoken) and "unknown" word, which is the remaining 20 keywords from the dataset. The dataset is split into training, validation and test set in the ratio of 80:10:10 while making sure that the audio clips from the same person stays in the same set. All models are trained in Google Tensorflow framework [31] using 5 Figure 4: Depthwise separable CNN architecture. the standard cross-entropy loss and Adam optimizer [32]. With a batch size of 100, the models are trained for 20K iterations with initial learning rate of $5 \times 10^{-4}$, and reduced to $10^{-4}$ after first 10K iterations. The training data is augmented with background noise and random time shift of up to 100ms. The trained models are evaluated based on the classification accuracy on the test set. 4.1 Training Results Table 2 summarizes the accuracy, memory requirement and operations per inference for the network architectures for KWS from literature [5, 6, 7, 8] trained on Google speech commands dataset [9]. For all the models, we use 40 MFCC features extracted from a speech frame of length 40ms with a stride of 20ms, which gives 1960 (49×40) features for 1 second of audio. The accuracy shown in the table is the accuracy on test set. The memory shown in the table assumes 8-bit weights and activations, which is sufficient to achieve same accuracy as that from a full-precision network. NN Architecture Accuracy Memory Operations DNN [5] 84.3% 288 KB 0.57 MOps CNN-1 [6] 90.7% 556 KB 76.02 MOps CNN-2 [6] 84.6% 149 KB 1.46 MOps LSTM [8] 88.8% 26 KB 2.06 MOps CRNN [7] 87.8% 298 KB 5.85 MOps Table 2: Neural network model accuracy. CNN-1, CNN-2 are (cnn-trad-fpool3, cnn-one-fstride4) architectures from [6]. Also, we assume that the memory for activations is reused across different layers and hence memory requirement for the activations uses the maximum of two consecutive layers. The operations in the table counts the total number of multiplications and additions in the matrix-multiplication operations in each layer in the network, which is representative of the execution time of the entire network. The models from the existing literature are optimized for different datasets and use different memory/compute resources, hence a direct comparison of accuracy is unfair. That said, these results still provide useful insights on the different neural network architectures for KWS: • Although DNNs do not achieve the best accuracy and tend to be memory intensive, they have less number of operations/inference and hence suit well to systems that have limited compute capability (e.g. systems running at low operating frequencies for energy-efficiency). • CNNs, on the other hand, achieve higher accuracy

than DNNs but at the cost of large number of operations and/or memory requirement. 6 • LSTMs and CRNNs achieve a balance between memory and operations while still achieving good accuracy.

4.2 Classifying Neural Networks for KWS based on Resource Requirements As discussed in section 2.2, memory footprint and execution time are the two important considerations in being able to run keyword spotting on microcontrollers. These should be considered when designing and optimizing neural networks for running keyword spotting. Based on typical microcontroller system configurations (as described in Table 1), we derive three sets of constraints for the neural networks in Table 3, targeting small, medium and large microcontroller systems. Both memory and compute limit are derived with assumptions that some amount of resources will be allocated for running other tasks such as OS, I/O, network communication, etc. The operations per inference limit assumes that the system is running 10 inferences per second. NN size NN memory limit Ops/inference limit Small (S) 80 KB 6 MOps Medium (M) 200 KB 20 MOps Large (L) 500 KB 80 MOps Table 3: Neural network (NN) classes for KWS models considered in this work, assuming 10 inferences per second and 8-bit weights/activations. 4.3 Resource Constrained Neural Network Architecture Exploration Figure 5 shows the number of operations per inference, memory requirement and test accuracy of neural network models from prior work [5, 6, 7, 8] trained on Google speech commands dataset overlayed with the memory and compute bounding boxes for the neural network classes from section 4.2. An ideal model would have high accuracy, small memory footprint and lower number of computations, i.e., close to the origin in Fig. 5. Apart from the LSTM model, the other models are too memory/compute resource heavy and do not fit into the bounding box S with 80KB/6MOps memory/compute limits. CNN-2, CRNN and DNN models fit in the M and L bounding boxes, but have lower accuracies as compared to the CNN-1 model, which does not fit in any of the boxes at all. The rest of this section discusses different hyperparameters of the feature extraction and neural network architectures that can be tuned in order to bring the models close to the origin and still achieve high accuracy. DNN:84.3% CNN-1:90.7% LSTM:88.8% CNN-2:84.6% CRNN:87.8% 100K 1M 10M 100M 10 100 1000 Operations Memory (KB) S(80KB, 6MOps) M(200KB, 20MOps) L(500KB, 80MOps) Figure 5: Number of operations vs. memory vs. test accuracy of NN models from prior work [5, 6, 7, 8] trained on the speech commands dataset [9]. As shown in Fig. 1, from each input speech signal, $T \times F$ features are extracted and the number of these features impact the model size, number of operations and accuracy. The key parameters in the feature extraction step that impact the model size, number of operations and accuracy are (1) number of MFCC features per frame (F) and (2) the frame stride (S). The number of MFCC features per audio frame (F) impacts the number of weights in fully-connected and recurrent layers, 7 but not in convolution layers as weights are reused in convolution layers. The frame stride (S), which determines the number of frames to be processed per inference (i.e. T), impacts the number of weights in fully-connected layers but not in recurrent and convolution layers because of the weight reuse. Both F and S impact the number of operations per inference. An efficient model would maximize accuracy using small $T \times F$, i.e., small F and/or large S. The neural network architectures and the corresponding hyperparameters explored in this work are summarized in Table 4. The LSTM model mentioned in the table includes peephole connections and output projection layer similar to that in [8], whereas basic LSTM model does not include those. CRNN uses one convolution layer followed by multi-layer GRU for the recurrent layers. We also use batch normalization for convolutional/fully-connected layers and layer normalization for recurrent layers. During inference, the parameters of batch normalization and layer normalization can be folded into the weights of the convolution or recurrent layers and hence these layers are ignored in memory/Ops computation. NN model Model hyperparameters DNN Number of fully-connected (FC) layers and size of each FC layer CNN Number of Conv layers: features/kernel size/stride, linear layer dim., FC layer size Basic LSTM Number of memory cells LSTM Number of memory cells, projection layer size GRU Number of memory cells CRNN Conv features/kernel size/stride, Number of GRU and

memory cells, FC layer size DS-CNN Number of DS-Conv layers, DS-Conv features/kernel size/stride Table 4: Neural network hyperparameters used in this study. We iteratively perform exhaustive search of feature extraction hyperparameters and NN model hyperparameters followed by manual selection to narrow down the search space. The final best performing models for each neural network architecture along with their memory requirements and operations are summarized in Table 5 and Fig. 6. The hyperparameters of these networks are summarized in Appendix A. From the results we can see that DNNs are memory-bound and achieve less accuracies and saturate at ~87% even when the model is scaled up. CNNs achieve better accuracies than DNN, but are limited by the weights in the final fully-connected layers. RNN models (i.e. Basic LSTM, LSTM and GRU) achieve better accuracies than CNNs and yield even smaller models with less Ops in some cases, demonstrating that exploiting temporal dependencies maximizes accuracy within the same resource budget. CRNN models, which combine the best properties of CNNs and RNNs, achieve better accuracies than both CNNs and RNNs, even with less Ops. CRNN architecture also scales up well when more memory/compute resources are available. DS-CNN achieves the best accuracies and demonstrate good scalability owing to their deeper architecture enabled by depthwise separable convolution layers, which are less compute/memory intensive. NN model S(80KB, 6MOps) M(200KB, 20MOps) L(500KB, 80MOps) Acc. Mem. Ops Acc. Mem. Ops Acc. Mem. Ops DNN 84.6% 80.0KB 158.8K 86.4% 199.4KB 397.0K 86.7% 496.6KB 990.2K CNN 91.6% 79.0KB 5.0M 92.2% 199.4KB 17.3M 92.7% 497.8KB 25.3M Basic LSTM 92.0% 63.3KB 5.9M 93.0% 196.5KB 18.9M 93.4% 494.5KB 47.9M LSTM 92.9% 79.5KB 3.9M 93.9% 198.6KB 19.2M 94.8% 498.8KB 48.4M GRU 93.5% 78.8KB 3.8M 94.2% 200.0KB 19.2M 94.7% 499.7KB 48.4M CRNN 94.0% 79.7KB 3.0M 94.4% 199.8KB 7.6M 95.0% 499.5KB 19.3M DS-CNN 94.4% 38.6KB 5.4M 94.9% 189.2KB 19.8M 95.4% 497.6KB 56.9M Table 5: Summary of best neural networks from the hyperparameter search. The memory required for storing the 8-bit weights and activations is shown in the table. To study the scalability of the models for smaller microcontroller systems with memory as low as 8KB, we expand the search space for DS-CNN models. Figure 7 shows the accuracy, memory/Ops requirements of the DS-CNN models targeted for such constrained devices. It shows that scaled-down DS-CNN models achieve better accuracies than DNN models with similar number of Ops, but with >10x reduction in memory requirement. 8 DNN:84.58% CNN:91.56% Basic LSTM:92.02% GRU:93.54% CRNN:94.09% DS-CNN:94.45% DNN:86.37% CNN:92.18% GRU:94.19% CRNN:94.38% DS-CNN:94.94% DNN:86.66% CNN:92.7% LSTM:94.81% CRNN:95% DS-CNN:95.39% 100K 1M 10M 100M 10 100 1000 Operations Memory (KB) M(200KB, 20MOps) L(500KB, 80MOps) S(80KB, 6MOps) Figure 6: Memory vs. Ops/inference of the best models described in Table 5. 75 80 85 90 95 100 4 16 64 256 Accuracy (%) Memory (KB) Accuracy (%) S(38.6KB, 5.4MOps) M(189.2KB, 19.8MOps) 256 75 80 85 90 95 100 200K 400K 800K 1.6M 3.2M 6.4M 12.8M 25.6M Accuracy (%) Operations Figure 7: Accuracy vs. memory and Ops of different DS-CNN models demonstrating the scalability of DS-CNN models down to <8KB memory footprint and <500K operations. 4.4 Neural Network Quantization Neural networks are typically trained with floating point weights and activations. Previous research [33, 34, 35] have shown that fixed-point weights is sufficient to run neural networks with minimal loss in accuracy. Microcontroller systems have limited memory, which motivates the quantization of 32-bit floating point weights to 8-bit fixed point weights for deployment, thus reducing the model size by 4×. Moreover, fixed-point integer operations run much faster than floating point operations in typical microcontrollers, which is another reason for executing quantized model during deployment. In this work, we use the quantization flow described in [34] using 8-bits to represent all the weights and activations. For a given signed 2's complement 8-bit fixed-point number, its value (v) can be expressed as $v = -B_7 . 2^{7-N} + \sum_{i=0}^{6} B_i . 2^{i-N}$, where N is the fractional length, which can also be negative. N is fixed for a given layer, but can be different in other layers. For example, N = 0 can represent the range [−128, 127] with a step of 1, N = 7 can represent the range [−1, 1 − (1/2^7)] with a step of 1/2^7 and N = −2

can represent the range [−512, 508] with a step of 2 2 . The weights are quantized to 8-bits progressively one layer at a time by finding the optimal N for each layer that minimizes the loss in accuracy because of quantization. After all the weights are quantized, the activations are also quantized in a similar way to find the appropriate fractional length N for each layer. Table 6 shows the accuracies of representative 8-bit networks quantized using this method and compared with those of the original full-precision networks. The table shows that the 9 accuracy of the quantized network is either same or marginally better than the full-precision network, possibly due to better regularization because of quantization. We believe that the same conclusion will hold for the other neural network models explored in this work. NN model 32-bit floating point model accuracy 8-bit quantized model accuracy Train Val. Test Train Val. Test DNN 97.77% 88.04% 86.66% 97.99% 88.91% 87.60% Basic LSTM 98.38% 92.69% 93.41% 98.21% 92.53% 93.51% GRU 99.23% 93.92% 94.68% 99.21% 93.66% 94.68% CRNN 98.34% 93.99% 95.00% 98.43% 94.08% 95.03% Table 6: Accuracy comparison of representative 8-bit quantized networks with full-precision networks. 4.5 KWS Deployment on Microcontroller We deployed the KWS application on Cortex-M7 based STM32F746G-DISCO development board using CMSIS-NN kernels [36]. A picture of the board performing KWS is shown in Fig. 8. The deployed model is a DNN model with 8-bit weights and 8-bit activations and KWS is running at 10 inferences per second. Each inference, including memory copying, MFCC feature extraction and DNN execution, takes about 12 ms. The microcontroller can be put into Wait-for-Interrupt (WFI) mode for the remaining time for power saving. The entire KWS application occupies ~70 KB memory, including ~66 KB for weights, ~1 KB for activations and ~ 2 KB for audio I/O and MFCC features." Visualize the state after each reasoning step.

ChatGPT

## Understanding the Error and Requirements

The error "Resizing cell port cnn_kws_accel.conv1.data_in from 16 bits to 320 bits" indicates a mismatch between the width of **data_in** port of **conv1** and the signal driving it (**mfcc_out**). Given the MFCC outputs 16-bit features, we need to ensure that the convolutional layer **conv1** correctly processes these features.

### Recalculating Port and Signal Widths

### MFCC Output

- MFCC output: 16-bit features (**MFCC_FEATURES** features, each **ACTIV_BITS** bits wide).

### Convolutional Layer Input

- **conv1.data_in** should be **MFCC_FEATURES * INPUT_CHANNELS * INPUT_HEIGHT * ACTIV_BITS**.

### Step-by-Step Solution

### Verify Parameters

- **MFCC_FEATURES** is 40.

- **INPUT_CHANNELS** is 1.

- **INPUT_HEIGHT** is 1.

- **ACTIV_BITS** is 16 (MFCC outputs 16-bit features).

Hence, **conv1.data_in** should be **40 \* 1 \* 1 \* 16 = 640** bits wide.

**Update the Top Module and MFCC Module**

1. **MFCC Module**: Ensure it outputs 16-bit wide features.

2. **Top Module (cnn_kws_accel)**: Adjust the connection and expected width.

**Update mfcc_accel Module**

verilog

Copy code

module mfcc_accel ( input wire clk, input wire rst_n, input wire [15:0] audio_sample, input wire sample_valid, output reg [MFCC_FEATURES*ACTIV_BITS-1:0] mfcc_feature, // Adjusted width output reg mfcc_valid ); // ... (implementation of MFCC feature extraction) endmodule

**Updated Top Module (cnn_kws_accel)**

verilog

Copy code

module cnn_kws_accel #( parameter NUM_KEYWORDS = 10, parameter MFCC_FEATURES = 40, parameter ACTIV_BITS = 16, // Updated to match MFCC output bit width parameter FC1_INPUT_SIZE = (MFCC_FEATURES/2)*(CONV2_NUM_FILTERS), parameter FC1_OUTPUT_SIZE = 64, parameter FC2_INPUT_SIZE = 64, parameter FC2_OUTPUT_SIZE = NUM_KEYWORDS, parameter CONV1_KERNEL_SIZE = 3, parameter CONV1_NUM_FILTERS = 8, parameter CONV2_KERNEL_SIZE = 3, parameter CONV2_NUM_FILTERS = 8 )( input wire clk, input wire rst_n, input wire [15:0] audio_in, input wire audio_valid, output reg [NUM_KEYWORDS-1:0] kws_result, output reg kws_valid, input wire [7:0] frame_size, input wire [7:8] frame_overlap, input wire [7:0] num_mfcc_coeffs, input wire [7:0] goertzel_coefs, output wire goertzel_coefs_start, output wire goertzel_coefs_valid, output wire goertzel_coefs_done, input wire [CONV1_NUM_FILTERS*CONV1_KERNEL_SIZE*CONV1_KERNEL_SIZE*ACTIV_BITS-1:0] conv1_weights, input wire [CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_biases, input wire conv1_load_weights, input wire conv1_load_biases, input wire [CONV2_NUM_FILTERS*CONV1_NUM_FILTERS*CONV2_KERNEL_SIZE*CONV2_KERNEL_SIZE*ACTIV_BITS-1:0] conv2_weights, input wire [CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_biases, input wire conv2_load_weights, input wire conv2_load_biases, input wire [FC1_OUTPUT_SIZE*FC1_INPUT_SIZE*ACTIV_BITS-1:0] fc1_weights, input wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_biases, input wire fc1_load_weights, input wire fc1_load_biases, input wire [FC2_OUTPUT_SIZE*FC2_INPUT_SIZE*ACTIV_BITS-1:0] fc2_weights, input wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_biases, input wire fc2_load_weights, input wire fc2_load_biases ); // MFCC module signals wire [MFCC_FEATURES*ACTIV_BITS-1:0] mfcc_out; // Adjusted width wire mfcc_valid; // CNN-KWS layers wire [MFCC_FEATURES*CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_out; wire conv1_valid; wire [MFCC_FEATURES*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_out; wire conv2_valid; wire [(MFCC_FEATURES/2)*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] maxpool_out; wire maxpool_valid; wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_out; wire fc1_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_out; wire fc2_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] softmax_out; wire softmax_valid; // MFCC module instantiation mfcc_accel mfcc ( .clk(clk), .rst(rst_n), .audio_sample(audio_in), .sample_valid(audio_valid), .mfcc_feature(mfcc_out),

.mfcc_valid(mfcc_valid) ); // Convolutional layer 1 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(1), .KERNEL_SIZE(CONV1_KERNEL_SIZE), .NUM_FILTERS(CONV1_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv1 ( .clk(clk), .rst_n(rst_n), .data_in(mfcc_out), // Corrected width .data_valid(mfcc_valid), .data_out(conv1_out), .data_out_valid(conv1_valid), .weights_in(conv1_weights), .biases_in(conv1_biases), .load_weights(conv1_load_weights), .load_biases(conv1_load_biases) ); // Convolutional layer 2 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV1_NUM_FILTERS), .KERNEL_SIZE(CONV2_KERNEL_SIZE), .NUM_FILTERS(CONV2_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv2 ( .clk(clk), .rst_n(rst_n), .data_in(conv1_out), .data_valid(conv1_valid), .data_out(conv2_out), .data_out_valid(conv2_valid), .weights_in(conv2_weights), .biases_in(conv2_biases), .load_weights(conv2_load_weights), .load_biases(conv2_load_biases) ); // Maxpooling layer maxpool2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV2_NUM_FILTERS), .KERNEL_SIZE(2), .STRIDE(2), .ACTIV_BITS(ACTIV_BITS) ) maxpool ( .clk(clk), .rst_n(rst_n), .data_in(conv2_out), .data_valid(conv2_valid), .data_out(maxpool_out), .data_out_valid(maxpool_valid) ); // Fully connected layer 1 fully_connected #( .INPUT_SIZE(FC1_INPUT_SIZE), .OUTPUT_SIZE(FC1_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc1 ( .clk(clk), .rst_n(rst_n), .data_in(maxpool_out), .data_valid(maxpool_valid), .data_out(fc1_out), .data_out_valid(fc1_valid), .weights_in(fc1_weights), .biases_in(fc1_biases), .load_weights(fc1_load_weights), .load_biases(fc1_load_biases ); // Fully connected layer 2 (output layer) fully_connected #( .INPUT_SIZE(FC2_INPUT_SIZE), .OUTPUT_SIZE(FC2_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc2 ( .clk(clk), .rst_n(rst_n), .data_in(fc1_out), .data_valid(fc1_valid), .data_out(fc2_out), .data_out_valid(fc2_valid), .weights_in(fc2_weights), .biases_in(fc2_biases), .load_weights(fc2_load_weights), .load_biases(fc2_load_biases ); // Softmax layer softmax #( .INPUT_SIZE(NUM_KEYWORDS), .ACTIV_BITS(ACTIV_BITS) ) softmax ( .clk(clk), .rst_n(rst_n), .data_in(fc2_out), .data_valid(fc2_valid), .data_out(softmax_out), .data_out_valid(softmax_valid) ); // Output assignment always @(posedge clk or negedge rst_n) begin if (!rst_n) begin kws_result <= 'b0; kws_valid <= 1'b0; end else begin kws_result <= softmax_out[NUM_KEYWORDS*ACTIV_BITS-1:0]; kws_valid <= softmax_valid; end end endmodule

## Summary

- The MFCC output width has been adjusted to 16 bits per feature, ensuring the total width is **MFCC_FEATURES * ACTIV_BITS = 40 * 16 = 640** bits.

- The convolutional layer **conv1**'s **data_in** port now correctly receives a 640-bit wide signal.

- Adjust the rest of the modules to handle the updated widths accordingly.

## Visualization and Verification

1. **Simulation**: Use a simulation tool to verify the signal widths and correctness of each module's output.

2. **Testbench**: Create a comprehensive testbench to validate the entire pipeline, ensuring all components work together seamlessly.

How did you come up with this "MFCC_FEATURES * ACTIV_BITS = 40 * 16 = 640 bits."?

ChatGPT

I apologize for the confusion in my previous response. The correct approach should account for the 16-bit integer MFCC outputs and how they fit into the convolutional layer input requirements. Let's go through the calculation again carefully.

**Understanding the Design Requirements**

1. **MFCC Output**: Each MFCC output is a 16-bit integer.

2. **CNN Input Requirements**:

    - **MFCC_FEATURES** is 40.

    - **INPUT_CHANNELS** is 1.

    - **INPUT_HEIGHT** is 1.

    - **ACTIV_BITS** is 16 (each MFCC output is 16 bits).

Given the parameters:

- **Width of MFCC output** = 16 bits (for each feature)

- **Number of MFCC features** = 40

So the **data_in** width for **conv1** should be **MFCC_FEATURES * ACTIV_BITS**.

**Calculations**

- **MFCC_FEATURES** = 40

- **ACTIV_BITS** = 16

So, **data_in** width for **conv1** should be:
data_in width=MFCC_FEATURES×ACTIV_BITSdata_in width=MFCC_FEATURES×ACTIV_BITS
data_in width=40×16=640 bitsdata_in width=40×16=640 bits

Let's update the code accordingly and ensure everything matches these requirements.

**Update the Top Module and MFCC Module**

**Ensure MFCC Module Outputs Correct Width**

verilog

Copy code

module mfcc_accel ( input wire clk, input wire rst_n, input wire [15:0] audio_sample, input wire sample_valid, output reg [MFCC_FEATURES*ACTIV_BITS-1:0] mfcc_feature, // Adjusted width output reg mfcc_valid ); // ... (implementation of MFCC feature extraction) endmodule

**Updated Top Module (cnn_kws_accel)**

verilog

Copy code

module cnn_kws_accel #( parameter NUM_KEYWORDS = 10, parameter MFCC_FEATURES = 40, parameter ACTIV_BITS = 16, // Updated to match MFCC output bit width parameter FC1_INPUT_SIZE = (MFCC_FEATURES/2)*(CONV2_NUM_FILTERS), parameter FC1_OUTPUT_SIZE = 64, parameter

FC2_INPUT_SIZE = 64, parameter FC2_OUTPUT_SIZE = NUM_KEYWORDS, parameter CONV1_KERNEL_SIZE = 3, parameter CONV1_NUM_FILTERS = 8, parameter CONV2_KERNEL_SIZE = 3, parameter CONV2_NUM_FILTERS = 8 )( input wire clk, input wire rst_n, input wire [15:0] audio_in, input wire audio_valid, output reg [NUM_KEYWORDS-1:0] kws_result, output reg kws_valid, input wire [7:0] frame_size, input wire [7:0] frame_overlap, input wire [7:0] num_mfcc_coeffs, input wire [7:0] goertzel_coefs, output wire goertzel_coefs_start, output wire goertzel_coefs_valid, output wire goertzel_coefs_done, input wire [CONV1_NUM_FILTERS*CONV1_KERNEL_SIZE*CONV1_KERNEL_SIZE*ACTIV_BITS-1:0] conv1_weights, input wire [CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_biases, input wire conv1_load_weights, input wire conv1_load_biases, input wire [CONV2_NUM_FILTERS*CONV1_NUM_FILTERS*CONV2_KERNEL_SIZE*CONV2_KERNEL_SIZE*ACTIV_BITS-1:0] conv2_weights, input wire [CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_biases, input wire conv2_load_weights, input wire conv2_load_biases, input wire [FC1_OUTPUT_SIZE*FC1_INPUT_SIZE*ACTIV_BITS-1:0] fc1_weights, input wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_biases, input wire fc1_load_weights, input wire fc1_load_biases, input wire [FC2_OUTPUT_SIZE*FC2_INPUT_SIZE*ACTIV_BITS-1:0] fc2_weights, input wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_biases, input wire fc2_load_weights, input wire fc2_load_biases ); // MFCC module signals wire [MFCC_FEATURES*ACTIV_BITS-1:0] mfcc_out; // Adjusted width wire mfcc_valid; // CNN-KWS layers wire [MFCC_FEATURES*CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_out; wire conv1_valid; wire [MFCC_FEATURES*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_out; wire conv2_valid; wire [(MFCC_FEATURES/2)*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] maxpool_out; wire maxpool_valid; wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_out; wire fc1_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_out; wire fc2_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] softmax_out; wire softmax_valid; // MFCC module instantiation mfcc_accel mfcc ( .clk(clk), .rst(rst_n), .audio_sample(audio_in), .sample_valid(audio_valid), .mfcc_feature(mfcc_out), .mfcc_valid(mfcc_valid) ); // Convolutional layer 1 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(1), .KERNEL_SIZE(CONV1_KERNEL_SIZE), .NUM_FILTERS(CONV1_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv1 ( .clk(clk), .rst_n(rst_n), .data_in(mfcc_out), // Corrected width .data_valid(mfcc_valid), .data_out(conv1_out), .data_out_valid(conv1_valid), .weights_in(conv1_weights), .biases_in(conv1_biases), .load_weights(conv1_load_weights), .load_biases(conv1_load_biases) ); // Convolutional layer 2 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV1_NUM_FILTERS), .KERNEL_SIZE(CONV2_KERNEL_SIZE), .NUM_FILTERS(CONV2_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv2 ( .clk(clk), .rst_n(rst_n), .data_in(conv1_out), .data_valid(conv1_valid), .data_out(conv2_out), .data_out_valid(conv2_valid), .weights_in(conv2_weights), .biases_in(conv2_biases), .load_weights(conv2_load_weights), .load_biases(conv2_load_biases) ); // Maxpooling layer maxpool2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV2_NUM_FILTERS), .KERNEL_SIZE(2), .STRIDE(2), .ACTIV_BITS(ACTIV_BITS) ) maxpool ( .clk(clk), .rst_n(rst_n), .data_in(conv2_out), .data_valid(conv2_valid), .data_out(maxpool_out), .data_out_valid(maxpool_valid) ); // Fully connected layer 1 fully_connected #( .INPUT_SIZE(FC1_INPUT_SIZE), .OUTPUT_SIZE(FC1_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc1 ( .clk(clk), .rst_n(rst_n), .data_in(maxpool_out), .data_valid(maxpool_valid), .data_out(fc1_out), .data_out_valid(fc1_valid), .weights_in(fc1_weights), .biases_in(fc1_biases), .load_weights(fc1_load_weights), .load_biases(fc1_load_biases ); // Fully connected layer 2 (output layer) fully_connected #( .INPUT_SIZE(FC2_INPUT_SIZE), .OUTPUT_SIZE(FC2_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc2 ( .clk(clk), .rst_n(rst_n),

.data_in(fc1_out), .data_valid(fc1_valid), .data_out(fc2_out), .data_out_valid(fc2_valid), .weights_in(fc2_weights), .biases_in(fc2_biases), .load_weights(fc2_load_weights), .load_biases(fc2_load_biases ); // Softmax layer softmax #( .INPUT_SIZE(NUM_KEYWORDS), .ACTIV_BITS(ACTIV_BITS) ) softmax ( .clk(clk), .rst_n(rst_n), .data_in(fc2_out), .data_valid(fc2_valid), .data_out(softmax_out), .data_out_valid(softmax_valid) ); // Output assignment always @(posedge clk or negedge rst_n) begin if (!rst_n) begin kws_result <= 'b0; kws_valid <= 1'b0; end else begin kws_result <= softmax_out[NUM_KEYWORDS*ACTIV_BITS-1:0]; kws_valid <= softmax_valid; end end endmodule

**Summary**

- Correctly adjusted the **mfcc_out** width to match the 16-bit feature requirement.

- Ensured that the **data_in** for **conv1** correctly receives 640 bits as input.

- Verified the calculation based on the design parameters and updated the modules accordingly.

By doing this, the design should now correctly handle the data width requirements, and the error should be resolved.

Compilation failed with : "Warning: Resizing cell port cnn_kws_accel.mfcc.mfcc_feature from 640 bits to 16 bits." mfcc feature is 16 bits. Refactor the rest port and signal widths based on this number.

ChatGPT

Thank you for the clarification. Since the MFCC feature output is 16 bits, we need to correctly handle the width of the data signals across the entire design to match this.

**Refactor Based on 16-bit MFCC Output**

Given that each MFCC feature output is 16 bits, we'll adjust the design so that the convolution layers and subsequent layers correctly interpret this 16-bit feature.

**Updated Top Module (cnn_kws_accel)**

verilog

Copy code

module cnn_kws_accel #( parameter NUM_KEYWORDS = 10, parameter MFCC_FEATURES = 40, // Number of MFCC features parameter ACTIV_BITS = 16, // Each MFCC feature is 16 bits parameter FC1_INPUT_SIZE = (MFCC_FEATURES/2)*(CONV2_NUM_FILTERS), parameter FC1_OUTPUT_SIZE = 64, parameter FC2_INPUT_SIZE = 64, parameter FC2_OUTPUT_SIZE = NUM_KEYWORDS, parameter CONV1_KERNEL_SIZE = 3, parameter CONV1_NUM_FILTERS = 8, parameter CONV2_KERNEL_SIZE = 3, parameter CONV2_NUM_FILTERS = 8 )( input wire clk, input wire rst_n, input wire [15:0] audio_in, input wire audio_valid, output reg [NUM_KEYWORDS-1:0] kws_result, output reg kws_valid, input wire [7:0] frame_size, input wire [7:0] frame_overlap, input wire [7:0] num_mfcc_coeffs, input wire [7:0] goertzel_coefs, output wire goertzel_coefs_start, output wire goertzel_coefs_valid, output wire goertzel_coefs_done, input wire [CONV1_NUM_FILTERS*CONV1_KERNEL_SIZE*CONV1_KERNEL_SIZE*ACTIV_BITS-1:0] conv1_weights, input wire [CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_biases, input wire conv1_load_weights, input wire conv1_load_biases, input wire [CONV2_NUM_FILTERS*CONV1_NUM_FILTERS*CONV2_KERNEL_SIZE*CONV2_KERNEL_SIZE*ACTIV_

BITS-1:0] conv2_weights, input wire [CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_biases, input wire conv2_load_weights, input wire conv2_load_biases, input wire [FC1_OUTPUT_SIZE*FC1_INPUT_SIZE*ACTIV_BITS-1:0] fc1_weights, input wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_biases, input wire fc1_load_weights, input wire fc1_load_biases, input wire [FC2_OUTPUT_SIZE*FC2_INPUT_SIZE*ACTIV_BITS-1:0] fc2_weights, input wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_biases, input wire fc2_load_weights, input wire fc2_load_biases ); // MFCC module signals wire [15:0] mfcc_out; // Adjusted width to 16 bits wire mfcc_valid; // CNN-KWS layers wire [MFCC_FEATURES*CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_out; wire conv1_valid; wire [MFCC_FEATURES*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_out; wire conv2_valid; wire [(MFCC_FEATURES/2)*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] maxpool_out; wire maxpool_valid; wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_out; wire fc1_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_out; wire fc2_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] softmax_out; wire softmax_valid; // MFCC module instantiation mfcc_accel mfcc ( .clk(clk), .rst(rst_n), .audio_sample(audio_in), .sample_valid(audio_valid), .mfcc_feature(mfcc_out), .mfcc_valid(mfcc_valid) ); // Convolutional layer 1 conv2d #( .INPUT_WIDTH(1), // Single feature .INPUT_HEIGHT(MFCC_FEATURES), // 40 features .INPUT_CHANNELS(1), .KERNEL_SIZE(CONV1_KERNEL_SIZE), .NUM_FILTERS(CONV1_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv1 ( .clk(clk), .rst_n(rst_n), .data_in(mfcc_out), // Corrected width .data_valid(mfcc_valid), .data_out(conv1_out), .data_out_valid(conv1_valid), .weights_in(conv1_weights), .biases_in(conv1_biases), .load_weights(conv1_load_weights), .load_biases(conv1_load_biases) ); // Convolutional layer 2 conv2d #( .INPUT_WIDTH(1), .INPUT_HEIGHT(MFCC_FEATURES), .INPUT_CHANNELS(CONV1_NUM_FILTERS), .KERNEL_SIZE(CONV2_KERNEL_SIZE), .NUM_FILTERS(CONV2_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv2 ( .clk(clk), .rst_n(rst_n), .data_in(conv1_out), .data_valid(conv1_valid), .data_out(conv2_out), .data_out_valid(conv2_valid), .weights_in(conv2_weights), .biases_in(conv2_biases), .load_weights(conv2_load_weights), .load_biases(conv2_load_biases) ); // Maxpooling layer maxpool2d #( .INPUT_WIDTH(1), .INPUT_HEIGHT(MFCC_FEATURES), .INPUT_CHANNELS(CONV2_NUM_FILTERS), .KERNEL_SIZE(2), .STRIDE(2), .ACTIV_BITS(ACTIV_BITS) ) maxpool ( .clk(clk), .rst_n(rst_n), .data_in(conv2_out), .data_valid(conv2_valid), .data_out(maxpool_out), .data_out_valid(maxpool_valid) ); // Fully connected layer 1 fully_connected #( .INPUT_SIZE(FC1_INPUT_SIZE), .OUTPUT_SIZE(FC1_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc1 ( .clk(clk), .rst_n(rst_n), .data_in(maxpool_out), .data_valid(maxpool_valid), .data_out(fc1_out), .data_out_valid(fc1_valid), .weights_in(fc1_weights), .biases_in(fc1_biases), .load_weights(fc1_load_weights), .load_biases(fc1_load_biases ); // Fully connected layer 2 (output layer) fully_connected #( .INPUT_SIZE(FC2_INPUT_SIZE), .OUTPUT_SIZE(FC2_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc2 ( .clk(clk), .rst_n(rst_n), .data_in(fc1_out), .data_valid(fc1_valid), .data_out(fc2_out), .data_out_valid(fc2_valid), .weights_in(fc2_weights), .biases_in(fc2_biases), .load_weights(fc2_load_weights), .load_biases(fc2_load_biases ); // Softmax layer softmax #( .INPUT_SIZE(NUM_KEYWORDS), .ACTIV_BITS(ACTIV_BITS) ) softmax ( .clk(clk), .rst_n(rst_n), .data_in(fc2_out), .data_valid(fc2_valid), .data_out(softmax_out), .data_out_valid(softmax_valid) ); // Output assignment always @(posedge clk or negedge rst_n) begin if (!rst_n) begin kws_result <= 'b0; kws_valid <= 1'b0; end else begin kws_result <= softmax_out[NUM_KEYWORDS*ACTIV_BITS-1:0]; kws_valid <= softmax_valid; end end endmodule

**Ensure Proper Signal Width Handling**

Since the MFCC output is 16 bits wide, let's make sure the rest of the design correctly handles this width.

**conv2d Module Update**

Ensure the **conv2d** module handles the input properly.

verilog

Copy code

```
module conv2d #( parameter INPUT_WIDTH = 1, // Single feature width parameter INPUT_HEIGHT =
40, // 40 MFCC features parameter INPUT_CHANNELS = 1, // Single channel input parameter
KERNEL_SIZE = 3, parameter NUM_FILTERS = 8, parameter PADDING = 1, parameter ACTIV_BITS = 16
) ( input wire clk, input wire rst_n, input wire
[INPUT_WIDTH*INPUT_HEIGHT*INPUT_CHANNELS*ACTIV_BITS-1:0] data_in, input wire data_valid,
output reg [INPUT_WIDTH*INPUT_HEIGHT*NUM_FILTERS*ACTIV_BITS-1:0] data_out, output reg
data_out_valid, input wire
[NUM_FILTERS*INPUT_CHANNELS*KERNEL_SIZE*KERNEL_SIZE*ACTIV_BITS-1:0] weights_in, input
wire [NUM_FILTERS*ACTIV_BITS-1:0] biases_in, input wire load_weights, input wire load_biases ); //
Declare weights and biases reg [ACTIV_BITS-1:0] weights [0:NUM_FILTERS-1][0:INPUT_CHANNELS-
1][0:KERNEL_SIZE-1][0:KERNEL_SIZE-1]; reg [ACTIV_BITS-1:0] biases [0:NUM_FILTERS-1]; // Declare
internal signals reg [ACTIV_BITS-1:0] input_buffer [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-
1][0:INPUT_CHANNELS-1]; reg [2*ACTIV_BITS-1:0] conv_result [0:INPUT_HEIGHT-
1][0:INPUT_WIDTH-1][0:NUM_FILTERS-1]; reg [ACTIV_BITS-1:0] relu_result [0:INPUT_HEIGHT-
1][0:INPUT_WIDTH-1][0:NUM_FILTERS-1]; // Load weights and biases integer i_load, j_load, k_load,
l_load; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset weights and biases for
(i_load = 0; i_load < NUM_FILTERS; i_load = i_load + 1) begin for (j_load = 0; j_load <
INPUT_CHANNELS; j_load = j_load + 1) begin for (k_load = 0; k_load < KERNEL_SIZE; k_load = k_load
+ 1) begin for (l_load = 0; l_load < KERNEL_SIZE; l_load = l_load + 1) begin
weights[i_load][j_load][k_load][l_load] <= 0; end end end biases[i_load] <= 0; end end else begin //
Load weights when load_weights is asserted if (load_weights) begin for (i_load = 0; i_load <
NUM_FILTERS; i_load = i_load + 1) begin for (j_load = 0; j_load < INPUT_CHANNELS; j_load = j_load +
1) begin for (k_load = 0; k_load < KERNEL_SIZE; k_load = k_load + 1) begin for (l_load = 0; l_load <
KERNEL_SIZE; l_load = l_load + 1) begin weights[i_load][j_load][k_load][l_load] <=
weights_in[(i_load*INPUT_CHANNELS*KERNEL_SIZE*KERNEL_SIZE +
j_load*KERNEL_SIZE*KERNEL_SIZE + k_load*KERNEL_SIZE + l_load)*ACTIV_BITS +: ACTIV_BITS]; end
end end end end // Load biases when load_biases is asserted if (load_biases) begin for (i_load = 0;
i_load < NUM_FILTERS; i_load = i_load + 1) begin biases[i_load] <= biases_in[i_load*ACTIV_BITS +:
ACTIV_BITS]; end end end end // Convolution operation integer m_conv, n_conv, p_conv, q_conv,
i_conv, j_conv, k_conv; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset
internal signals and output for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for
(j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv <
INPUT_CHANNELS; k_conv = k_conv + 1) begin input_buffer[i_conv][j_conv][k_conv] <= 0; end for
(m_conv = 0; m_conv < NUM_FILTERS; m_conv = m_conv + 1) begin
conv_result[i_conv][j_conv][m_conv] <= 0; relu_result[i_conv][j_conv][m_conv] <= 0; end end end
data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Shift input data into buffer for
(i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv <
INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < INPUT_CHANNELS; k_conv =
k_conv + 1) begin if (j_conv < INPUT_WIDTH - 1) begin input_buffer[i_conv][j_conv][k_conv] <=
input_buffer[i_conv][j_conv+1][k_conv]; end else begin input_buffer[i_conv][j_conv][k_conv] <=
data_in[i_conv*INPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS +
j_conv*INPUT_CHANNELS*ACTIV_BITS + k_conv*ACTIV_BITS +: ACTIV_BITS]; end end end end //
```

Perform convolution for (m_conv = 0; m_conv < INPUT_HEIGHT; m_conv = m_conv + 1) begin for (n_conv = 0; n_conv < INPUT_WIDTH; n_conv = n_conv + 1) begin for (p_conv = 0; p_conv < NUM_FILTERS; p_conv = p_conv + 1) begin conv_result[m_conv][n_conv][p_conv] = {{(2*ACTIV_BITS-ACTIV_BITS){1'b0}}, biases[p_conv]}; for (q_conv = 0; q_conv < INPUT_CHANNELS; q_conv = q_conv + 1) begin for (i_conv = 0; i_conv < KERNEL_SIZE; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < KERNEL_SIZE; j_conv = j_conv + 1) begin if (m_conv + i_conv - PADDING >= 0 && m_conv + i_conv - PADDING < INPUT_HEIGHT && n_conv + j_conv - PADDING >= 0 && n_conv + j_conv - PADDING < INPUT_WIDTH) begin conv_result[m_conv][n_conv][p_conv] = conv_result[m_conv][n_conv][p_conv] + weights[p_conv][q_conv][i_conv][j_conv] * input_buffer[m_conv + i_conv - PADDING][n_conv + j_conv - PADDING][q_conv]; end end end end end end end // Apply ReLU activation for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < NUM_FILTERS; k_conv = k_conv + 1) begin relu_result[i_conv][j_conv][k_conv] = (conv_result[i_conv][j_conv][k_conv][2*ACTIV_BITS-1] == 0) ? conv_result[i_conv][j_conv][k_conv][ACTIV_BITS-1:0] : 0; end end end // Assign output for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < NUM_FILTERS; k_conv = k_conv + 1) begin data_out[i_conv*INPUT_WIDTH*NUM_FILTERS*ACTIV_BITS + j_conv*NUM_FILTERS*ACTIV_BITS + k_conv*ACTIV_BITS +: ACTIV_BITS] <= relu_result[i_conv][j_conv][k_conv]; end end end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule

**Ensure Proper Signal Width Handling**

verilog

Copy code

module fully_connected #( parameter INPUT_SIZE = 320, // Adjust based on new dimensions parameter OUTPUT_SIZE = 64, parameter ACTIV_BITS = 16 ) ( input wire clk, input wire rst_n, input wire [INPUT_SIZE*ACTIV_BITS-1:0] data_in, input wire data_valid, output reg [OUTPUT_SIZE*ACTIV_BITS-1:0] data_out, output reg data_out_valid, input wire [OUTPUT_SIZE*INPUT_SIZE*ACTIV_BITS-1:0] weights_in, input wire [OUTPUT_SIZE*ACTIV_BITS-1:0] biases_in, input wire load_weights, input wire load_biases ); // Declare weights and biases reg [ACTIV_BITS-1:0] weights [0:OUTPUT_SIZE-1][0:INPUT_SIZE-1]; reg [ACTIV_BITS-1:0] biases [0:OUTPUT_SIZE-1]; // Declare internal signals reg [2*ACTIV_BITS-1:0] acc_result [0:OUTPUT_SIZE-1]; reg [ACTIV_BITS-1:0] relu_result [0:OUTPUT_SIZE-1]; // Load weights and biases integer i_load, j_load; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset weights and biases for (i_load = 0; i_load < OUTPUT_SIZE; i_load = i_load + 1) begin for (j_load = 0; j_load < INPUT_SIZE; j_load = j_load + 1) begin weights[i_load][j_load] <= 0; end biases[i_load] <= 0; end end else begin // Load weights when load_weights is asserted if (load_weights) begin for (i_load = 0; i_load < OUTPUT_SIZE; i_load = i_load + 1) begin for (j_load = 0; j_load < INPUT_SIZE; j_load = j_load + 1) begin weights[i_load][j_load] <= weights_in[(i_load*INPUT_SIZE + j_load)*ACTIV_BITS +: ACTIV_BITS]; end end end // Load biases when load_biases is asserted if (load_biases) begin for (i_load = 0; i_load < OUTPUT_SIZE; i_load = i_load + 1) begin biases[i_load] <= biases_in[i_load*ACTIV_BITS +: ACTIV_BITS]; end end end end // Fully connected layer operation integer i_fc, j_fc; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset internal signals and output for (i_fc = 0; i_fc < OUTPUT_SIZE; i_fc = i_fc + 1) begin acc_result[i_fc] <= 0; relu_result[i_fc] <= 0; end data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Perform matrix multiplication for (i_fc = 0; i_fc < OUTPUT_SIZE; i_fc = i_fc + 1) begin acc_result[i_fc] = {{(2*ACTIV_BITS-ACTIV_BITS){1'b0}}, biases[i_fc]}; for (j_fc = 0; j_fc < INPUT_SIZE; j_fc = j_fc + 1)

begin acc_result[i_fc] = acc_result[i_fc] + weights[i_fc][j_fc] * data_in[j_fc*ACTIV_BITS +: ACTIV_BITS]; end end // Apply ReLU activation for (i_fc = 0; i_fc < OUTPUT_SIZE; i_fc = i_fc + 1) begin relu_result[i_fc] <= (acc_result[i_fc][2*ACTIV_BITS-1] == 0) ? acc_result[i_fc][ACTIV_BITS-1:0] : 0; end // Assign output for (i_fc = 0; i_fc < OUTPUT_SIZE; i_fc = i_fc + 1) begin data_out[i_fc*ACTIV_BITS +: ACTIV_BITS] <= relu_result[i_fc]; end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule

**Ensure Proper Signal Width Handling for maxpool2d**

verilog

Copy code

module maxpool2d #( parameter INPUT_WIDTH = 1, parameter INPUT_HEIGHT = 40, parameter INPUT_CHANNELS = 8, parameter KERNEL_SIZE = 2, parameter STRIDE = 2, parameter ACTIV_BITS = 16 ) ( input wire clk, input wire rst_n, input wire [INPUT_WIDTH*INPUT_HEIGHT*INPUT_CHANNELS*ACTIV_BITS-1:0] data_in, input wire data_valid, output reg [(INPUT_WIDTH/STRIDE)*INPUT_CHANNELS*ACTIV_BITS-1:0] data_out, output reg data_out_valid ); localparam OUTPUT_WIDTH = INPUT_WIDTH / STRIDE; localparam OUTPUT_HEIGHT = INPUT_HEIGHT / STRIDE; // Declare internal signals reg [ACTIV_BITS-1:0] input_buffer [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:INPUT_CHANNELS-1]; reg [ACTIV_BITS-1:0] max_value [0:INPUT_CHANNELS-1]; // Max pooling operation integer i, j, k, m, n; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset internal signals and output for (i = 0; i < INPUT_HEIGHT; i = i + 1) begin for (j = 0; j < INPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin input_buffer[i][j][k] <= 0; end end end data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Shift input data into buffer for (i = 0; i < INPUT_HEIGHT; i = i + 1) begin for (j = 0; j < INPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin if (j < INPUT_WIDTH - 1) begin input_buffer[i][j][k] <= input_buffer[i][j+1][k]; end else begin input_buffer[i][j][k] <= data_in[i*INPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS + j*INPUT_CHANNELS*ACTIV_BITS + k*ACTIV_BITS +: ACTIV_BITS]; end end end end // Perform max pooling for (i = 0; i < OUTPUT_HEIGHT; i = i + 1) begin for (j = 0; j < OUTPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin max_value[k] = input_buffer[i*STRIDE][j*STRIDE][k]; for (m = 0; m < KERNEL_SIZE; m = m + 1) begin for (n = 0; n < KERNEL_SIZE; n = n + 1) begin if (i*STRIDE + m < INPUT_HEIGHT && j*STRIDE + n < INPUT_WIDTH) begin max_value[k] = (input_buffer[i*STRIDE + m][j*STRIDE + n][k] > max_value[k]) ? input_buffer[i*STRIDE + m][j*STRIDE + n][k] : max_value[k]; end end end data_out[i*OUTPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS + j*INPUT_CHANNELS*ACTIV_BITS + k*ACTIV_BITS +: ACTIV_BITS] <= max_value[k]; end end end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule

**Conclusion**

By ensuring that each module handles 16-bit MFCC features correctly, the design should now compile without width mismatch errors. Each module was refactored to ensure proper signal width handling throughout the design.

Warning: Resizing cell port cnn_kws_accel.maxpool.data_out from 2560 bits to 2 bits. Warning: Resizing cell port cnn_kws_accel.conv1.data_in from 16 bits to 640 bits.

ChatGPT

I apologize for the confusion in my previous responses. Given that the MFCC output is 16 bits wide, we need to ensure that the width of **mfcc_out** correctly matches the expected input width for the convolutional layer. Let's refactor the widths and ensure the calculations are correct.

**Correcting the Widths**

1. **MFCC Output**: Each MFCC feature output is 16 bits.

2. **Convolutional Layer Input**: The convolutional layer should take all MFCC features as input. Given 40 features, and each feature being 16 bits, the total width should be 40×16=64040×16=640 bits.

**Refactor the Top Module and Submodules**

**Ensure Proper Signal Width Handling**

**Update Top Module (cnn_kws_accel)**

Adjust the signal widths and parameters accordingly:

verilog

Copy code

module cnn_kws_accel #( parameter NUM_KEYWORDS = 10, parameter MFCC_FEATURES = 40, parameter ACTIV_BITS = 16, // Each MFCC feature is 16 bits parameter FC1_INPUT_SIZE = (MFCC_FEATURES/2) * (CONV2_NUM_FILTERS), parameter FC1_OUTPUT_SIZE = 64, parameter FC2_INPUT_SIZE = 64, parameter FC2_OUTPUT_SIZE = NUM_KEYWORDS, parameter CONV1_KERNEL_SIZE = 3, parameter CONV1_NUM_FILTERS = 8, parameter CONV2_KERNEL_SIZE = 3, parameter CONV2_NUM_FILTERS = 8 )( input wire clk, input wire rst_n, input wire [15:0] audio_in, input wire audio_valid, output reg [NUM_KEYWORDS-1:0] kws_result, output reg kws_valid, input wire [7:0] frame_size, input wire [7:0] frame_overlap, input wire [7:0] num_mfcc_coeffs, input wire [7:0] goertzel_coefs, output wire goertzel_coefs_start, output wire goertzel_coefs_valid, output wire goertzel_coefs_done, input wire [CONV1_NUM_FILTERS * CONV1_KERNEL_SIZE * CONV1_KERNEL_SIZE * ACTIV_BITS-1:0] conv1_weights, input wire [CONV1_NUM_FILTERS * ACTIV_BITS-1:0] conv1_biases, input wire conv1_load_weights, input wire conv1_load_biases, input wire [CONV2_NUM_FILTERS * CONV1_NUM_FILTERS * CONV2_KERNEL_SIZE * CONV2_KERNEL_SIZE * ACTIV_BITS-1:0] conv2_weights, input wire [CONV2_NUM_FILTERS * ACTIV_BITS-1:0] conv2_biases, input wire conv2_load_weights, input wire conv2_load_biases, input wire [FC1_OUTPUT_SIZE * FC1_INPUT_SIZE * ACTIV_BITS-1:0] fc1_weights, input wire [FC1_OUTPUT_SIZE * ACTIV_BITS-1:0] fc1_biases, input wire fc1_load_weights, input wire fc1_load_biases, input wire [FC2_OUTPUT_SIZE * FC2_INPUT_SIZE * ACTIV_BITS-1:0] fc2_weights, input wire [FC2_OUTPUT_SIZE * ACTIV_BITS-1:0] fc2_biases, input wire fc2_load_weights, input wire fc2_load_biases ); // MFCC module signals wire [MFCC_FEATURES * ACTIV_BITS-1:0] mfcc_out; // 640 bits wire mfcc_valid; // CNN-KWS layers wire [(MFCC_FEATURES * CONV1_NUM_FILTERS * ACTIV_BITS)-1:0] conv1_out; wire conv1_valid; wire [(MFCC_FEATURES * CONV2_NUM_FILTERS * ACTIV_BITS)-1:0] conv2_out; wire conv2_valid; wire [((MFCC_FEATURES/2) * CONV2_NUM_FILTERS * ACTIV_BITS)-1:0] maxpool_out; wire maxpool_valid; wire [(FC1_OUTPUT_SIZE * ACTIV_BITS)-1:0] fc1_out; wire fc1_valid; wire [(FC2_OUTPUT_SIZE * ACTIV_BITS)-1:0] fc2_out; wire fc2_valid; wire [(FC2_OUTPUT_SIZE * ACTIV_BITS)-1:0] softmax_out; wire softmax_valid; // MFCC module instantiation mfcc_accel mfcc ( .clk(clk), .rst(rst_n), .audio_sample(audio_in), .sample_valid(audio_valid), .mfcc_feature(mfcc_out), // 640 bits .mfcc_valid(mfcc_valid) ); // Convolutional layer 1 conv2d #(

.INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(1), .KERNEL_SIZE(CONV1_KERNEL_SIZE), .NUM_FILTERS(CONV1_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv1 ( .clk(clk), .rst_n(rst_n), .data_in(mfcc_out), // 640 bits .data_valid(mfcc_valid), .data_out(conv1_out), .data_out_valid(conv1_valid), .weights_in(conv1_weights), .biases_in(conv1_biases), .load_weights(conv1_load_weights), .load_biases(conv1_load_biases) ); // Convolutional layer 2 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV1_NUM_FILTERS), .KERNEL_SIZE(CONV2_KERNEL_SIZE), .NUM_FILTERS(CONV2_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv2 ( .clk(clk), .rst_n(rst_n), .data_in(conv1_out), .data_valid(conv1_valid), .data_out(conv2_out), .data_out_valid(conv2_valid), .weights_in(conv2_weights), .biases_in(conv2_biases), .load_weights(conv2_load_weights), .load_biases(conv2_load_biases) ); // Maxpooling layer maxpool2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV2_NUM_FILTERS), .KERNEL_SIZE(2), .STRIDE(2), .ACTIV_BITS(ACTIV_BITS) ) maxpool ( .clk(clk), .rst_n(rst_n), .data_in(conv2_out), .data_valid(conv2_valid), .data_out(maxpool_out), .data_out_valid(maxpool_valid) ); // Fully connected layer 1 fully_connected #( .INPUT_SIZE(FC1_INPUT_SIZE), .OUTPUT_SIZE(FC1_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc1 ( .clk(clk), .rst_n(rst_n), .data_in(maxpool_out), .data_valid(maxpool_valid), .data_out(fc1_out), .data_out_valid(fc1_valid), .weights_in(fc1_weights), .biases_in(fc1_biases), .load_weights(fc1_load_weights), .load_biases(fc1_load_biases) ); // Fully connected layer 2 (output layer) fully_connected #( .INPUT_SIZE(FC2_INPUT_SIZE), .OUTPUT_SIZE(FC2_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc2 ( .clk(clk), .rst_n(rst_n), .data_in(fc1_out), .data_valid(fc1_valid), .data_out(fc2_out), .data_out_valid(fc2_valid), .weights_in(fc2_weights), .biases_in(fc2_biases), .load_weights(fc2_load_weights), .load_biases(fc2_load_biases) ); // Softmax layer softmax #( .INPUT_SIZE(NUM_KEYWORDS), .ACTIV_BITS(ACTIV_BITS) ) softmax ( .clk(clk), .rst_n(rst_n), .data_in(fc2_out), .data_valid(fc2_valid), .data_out(softmax_out), .data_out_valid(softmax_valid) ); // Output assignment always @(posedge clk or negedge rst_n) begin if (!rst_n) begin kws_result <= 'b0; kws_valid <= 1'b0; end else begin kws_result <= softmax_out[NUM_KEYWORDS*ACTIV_BITS-1:0]; kws_valid <= softmax_valid; end end endmodule

**Ensure Proper Signal Width Handling in Submodules**

**Update the conv2d Module**

verilog

Copy code

module conv2d #( parameter INPUT_WIDTH = 40, // MFCC features parameter INPUT_HEIGHT = 1, // Single feature height parameter INPUT_CHANNELS = 1, // Single channel input parameter KERNEL_SIZE = 3, parameter NUM_FILTERS = 8, parameter PADDING = 1, parameter ACTIV_BITS = 16 ) ( input wire clk, input wire rst_n, input wire [INPUT_WIDTH * INPUT_HEIGHT * INPUT_CHANNELS * ACTIV_BITS-1:0] data_in, input wire data_valid, output reg [INPUT_WIDTH * INPUT_HEIGHT * NUM_FILTERS * ACTIV_BITS-1:0] data_out, output reg data_out_valid, input wire [NUM_FILTERS * INPUT_CHANNELS * KERNEL_SIZE * KERNEL_SIZE * ACTIV_BITS-1:0] weights_in, input wire [NUM_FILTERS * ACTIV_BITS-1:0] biases_in, input wire load_weights, input wire load_biases ); // Declare weights and biases reg [ACTIV_BITS-1:0] weights [0:NUM_FILTERS-1][0:INPUT_CHANNELS-1][0:KERNEL_SIZE-1][0:KERNEL_SIZE-1]; reg [ACTIV_BITS-1:0] biases [0:NUM_FILTERS-1]; // Declare internal signals reg [ACTIV_BITS-1:0] input_buffer [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:INPUT_CHANNELS-1]; reg [2*ACTIV_BITS-1:0] conv_result [0:INPUT_HEIGHT-

```verilog
1][0:INPUT_WIDTH-1][0:NUM_FILTERS-1]; reg [ACTIV_BITS-1:0] relu_result [0:INPUT_HEIGHT-
1][0:INPUT_WIDTH-1][0:NUM_FILTERS-1]; // Load weights and biases integer i_load, j_load, k_load,
l_load; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset weights and biases for
(i_load = 0; i_load < NUM_FILTERS; i_load = i_load + 1) begin for (j_load = 0; j_load <
INPUT_CHANNELS; j_load = j_load + 1) begin for (k_load = 0; k_load < KERNEL_SIZE; k_load = k_load
+ 1) begin for (l_load = 0; l_load < KERNEL_SIZE; l_load = l_load + 1) begin
weights[i_load][j_load][k_load][l_load] <= 0; end end end biases[i_load] <= 0; end end else begin //
Load weights when load_weights is asserted if (load_weights) begin for (i_load = 0; i_load <
NUM_FILTERS; i_load = i_load + 1) begin for (j_load = 0; j_load < INPUT_CHANNELS; j_load = j_load +
1) begin for (k_load = 0; k_load < KERNEL_SIZE; k_load = k_load + 1) begin for (l_load = 0; l_load <
KERNEL_SIZE; l_load = l_load + 1) begin weights[i_load][j_load][k_load][l_load] <=
weights_in[(i_load*INPUT_CHANNELS*KERNEL_SIZE*KERNEL_SIZE +
j_load*KERNEL_SIZE*KERNEL_SIZE + k_load*KERNEL_SIZE + l_load)*ACTIV_BITS +: ACTIV_BITS]; end
end end end end // Load biases when load_biases is asserted if (load_biases) begin for (i_load = 0;
i_load < NUM_FILTERS; i_load = i_load + 1) begin biases[i_load] <= biases_in[i_load*ACTIV_BITS +:
ACTIV_BITS]; end end end end // Convolution operation integer m_conv, n_conv, p_conv, q_conv,
i_conv, j_conv, k_conv; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset
internal signals and output for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for
(j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv <
INPUT_CHANNELS; k_conv = k_conv + 1) begin input_buffer[i_conv][j_conv][k_conv] <= 0; end for
(m_conv = 0; m_conv < NUM_FILTERS; m_conv = m_conv + 1) begin
conv_result[i_conv][j_conv][m_conv] <= 0; relu_result[i_conv][j_conv][m_conv] <= 0; end end end
data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Shift input data into buffer for
(i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv <
INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < INPUT_CHANNELS; k_conv =
k_conv + 1) begin if (j_conv < INPUT_WIDTH - 1) begin input_buffer[i_conv][j_conv][k_conv] <=
input_buffer[i_conv][j_conv+1][k_conv]; end else begin input_buffer[i_conv][j_conv][k_conv] <=
data_in[i_conv*INPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS +
j_conv*INPUT_CHANNELS*ACTIV_BITS + k_conv*ACTIV_BITS +: ACTIV_BITS]; end end end end //
Perform convolution for (m_conv = 0; m_conv < INPUT_HEIGHT; m_conv = m_conv + 1) begin for
(n_conv = 0; n_conv < INPUT_WIDTH; n_conv = n_conv + 1) begin for (p_conv = 0; p_conv <
NUM_FILTERS; p_conv = p_conv + 1) begin conv_result[m_conv][n_conv][p_conv] = {{(2*ACTIV_BITS-
ACTIV_BITS){1'b0}}, biases[p_conv]}; for (q_conv = 0; q_conv < INPUT_CHANNELS; q_conv = q_conv +
1) begin for (i_conv = 0; i_conv < KERNEL_SIZE; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv <
KERNEL_SIZE; j_conv = j_conv + 1) begin if (m_conv + i_conv - PADDING >= 0 && m_conv + i_conv -
PADDING < INPUT_HEIGHT && n_conv + j_conv - PADDING >= 0 && n_conv + j_conv - PADDING <
INPUT_WIDTH) begin conv_result[m_conv][n_conv][p_conv] =
conv_result[m_conv][n_conv][p_conv] + weights[p_conv][q_conv][i_conv][j_conv] *
input_buffer[m_conv + i_conv - PADDING][n_conv + j_conv - PADDING][q_conv]; end end end end
end end end // Apply ReLU activation for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1)
begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv <
NUM_FILTERS; k_conv = k_conv + 1) begin relu_result[i_conv][j_conv][k_conv] =
(conv_result[i_conv][j_conv][k_conv][2*ACTIV_BITS-1] == 0) ?
conv_result[i_conv][j_conv][k_conv][ACTIV_BITS-1:0] : 0; end end end // Assign output for (i_conv =
0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv
= j_conv + 1) begin for (k_conv = 0; k_conv < NUM_FILTERS; k_conv = k_conv + 1) begin
data_out[i_conv*INPUT_WIDTH*NUM_FILTERS*ACTIV_BITS + j_conv*NUM_FILTERS*ACTIV_BITS +
```

k_conv*ACTIV_BITS +: ACTIV_BITS] <= relu_result[i_conv][j_conv][k_conv]; end end end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule

**Update the maxpool2d Module**

verilog

Copy code

```
module maxpool2d #( parameter INPUT_WIDTH = 40, parameter INPUT_HEIGHT = 1, parameter INPUT_CHANNELS = 8, parameter KERNEL_SIZE = 2, parameter STRIDE = 2, parameter ACTIV_BITS = 16 ) ( input wire clk, input wire rst_n, input wire [INPUT_WIDTH * INPUT_HEIGHT * INPUT_CHANNELS * ACTIV_BITS-1:0] data_in, input wire data_valid, output reg [(INPUT_WIDTH/STRIDE) * INPUT_CHANNELS * ACTIV_BITS-1:0] data_out, output reg data_out_valid ); localparam OUTPUT_WIDTH = INPUT_WIDTH / STRIDE; localparam OUTPUT_HEIGHT = INPUT_HEIGHT / STRIDE; // Declare internal signals reg [ACTIV_BITS-1:0] input_buffer [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:INPUT_CHANNELS-1]; reg [ACTIV_BITS-1:0] max_value [0:INPUT_CHANNELS-1]; // Max pooling operation integer i, j, k, m, n; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset internal signals and output for (i = 0; i < INPUT_HEIGHT; i = i + 1) begin for (j = 0; j < INPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin input_buffer[i][j][k] <= 0; end end end data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Shift input data into buffer for (i = 0; i < INPUT_HEIGHT; i = i + 1) begin for (j = 0; j < INPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin if (j < INPUT_WIDTH - 1) begin input_buffer[i][j][k] <= input_buffer[i][j+1][k]; end else begin input_buffer[i][j][k] <= data_in[i*INPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS + j*INPUT_CHANNELS*ACTIV_BITS + k*ACTIV_BITS +: ACTIV_BITS]; end end end end // Perform max pooling for (i = 0; i < OUTPUT_HEIGHT; i = i + 1) begin for (j = 0; j < OUTPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin max_value[k] = input_buffer[i*STRIDE][j*STRIDE][k]; for (m = 0; m < KERNEL_SIZE; m = m + 1) begin for (n = 0; n < KERNEL_SIZE; n = n + 1) begin if (i*STRIDE + m < INPUT_HEIGHT && j*STRIDE + n < INPUT_WIDTH) begin max_value[k] = (input_buffer[i*STRIDE + m][j*STRIDE + n][k] > max_value[k]) ? input_buffer[i*STRIDE + m][j*STRIDE + n][k] : max_value[k]; end end end data_out[i*OUTPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS + j*INPUT_CHANNELS*ACTIV_BITS + k*ACTIV_BITS +: ACTIV_BITS] <= max_value[k]; end end end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule
```

**Conclusion**

By ensuring that each module correctly handles the 16-bit MFCC features, the design should now compile without width mismatch errors. The adjustments in the **mfcc_accel**, **conv2d**, and **maxpool2d** modules ensure proper signal width handling throughout the design.

I get the following compilation error "Resizing cell port cnn_kws_accel.mfcc.mfcc_feature from 640 bits to 16 bits.". Visualize the state after each reasoning step.

ChatGPT

Given the error message "Resizing cell port cnn_kws_accel.mfcc.mfcc_feature from 640 bits to 16 bits", it's clear that there is a discrepancy between the expected and actual widths of the **mfcc_feature** signal. Let's revisit the design requirements and correct the signal widths.

**Understanding the Requirements**

1. **MFCC Output**:

   - Each MFCC feature is 16 bits wide.

   - There are 40 MFCC features.

2. **Signal Widths**:

   - The **mfcc_feature** signal should be **MFCC_FEATURES * ACTIV_BITS** wide.

**Refactor the Top Module and MFCC Module**

**Ensure Proper Signal Width Handling**

**Correct the MFCC Module**

Let's define the **mfcc_accel** module properly:

verilog

Copy code

module mfcc_accel ( input wire clk, input wire rst_n, input wire [15:0] audio_sample, input wire sample_valid, output reg [639:0] mfcc_feature, // 40 features * 16 bits = 640 bits output reg mfcc_valid ); // Implementation of MFCC feature extraction endmodule

**Update Top Module (cnn_kws_accel)**

Adjust the signal widths and parameters accordingly:

verilog

Copy code

module cnn_kws_accel #( parameter NUM_KEYWORDS = 10, parameter MFCC_FEATURES = 40, parameter ACTIV_BITS = 16, // Each MFCC feature is 16 bits parameter FC1_INPUT_SIZE = (MFCC_FEATURES/2) * (CONV2_NUM_FILTERS), parameter FC1_OUTPUT_SIZE = 64, parameter FC2_INPUT_SIZE = 64, parameter FC2_OUTPUT_SIZE = NUM_KEYWORDS, parameter CONV1_KERNEL_SIZE = 3, parameter CONV1_NUM_FILTERS = 8, parameter CONV2_KERNEL_SIZE = 3, parameter CONV2_NUM_FILTERS = 8 )( input wire clk, input wire rst_n, input wire [15:0] audio_in, input wire audio_valid, output reg [NUM_KEYWORDS-1:0] kws_result, output reg kws_valid, input wire [7:0] frame_size, input wire [7:0] frame_overlap, input wire [7:0] num_mfcc_coeffs, input wire [7:0] goertzel_coefs, output wire goertzel_coefs_start, output wire goertzel_coefs_valid, output wire goertzel_coefs_done, input wire [CONV1_NUM_FILTERS * CONV1_KERNEL_SIZE * CONV1_KERNEL_SIZE * ACTIV_BITS-1:0] conv1_weights, input wire [CONV1_NUM_FILTERS * ACTIV_BITS-1:0] conv1_biases, input wire conv1_load_weights, input wire conv1_load_biases, input wire [CONV2_NUM_FILTERS * CONV1_NUM_FILTERS * CONV2_KERNEL_SIZE * CONV2_KERNEL_SIZE * ACTIV_BITS-1:0] conv2_weights, input wire [CONV2_NUM_FILTERS * ACTIV_BITS-1:0] conv2_biases, input wire conv2_load_weights, input wire conv2_load_biases, input wire [FC1_OUTPUT_SIZE * FC1_INPUT_SIZE * ACTIV_BITS-1:0] fc1_weights, input wire [FC1_OUTPUT_SIZE * ACTIV_BITS-1:0] fc1_biases, input wire fc1_load_weights, input wire fc1_load_biases, input wire [FC2_OUTPUT_SIZE * FC2_INPUT_SIZE * ACTIV_BITS-1:0] fc2_weights, input wire [FC2_OUTPUT_SIZE * ACTIV_BITS-1:0] fc2_biases, input wire fc2_load_weights, input wire fc2_load_biases ); // MFCC module signals wire [MFCC_FEATURES * ACTIV_BITS-1:0] mfcc_out; // 640 bits wire mfcc_valid; // CNN-KWS layers wire [(MFCC_FEATURES * CONV1_NUM_FILTERS * ACTIV_BITS)-1:0] conv1_out; wire

conv1_valid; wire [(MFCC_FEATURES * CONV2_NUM_FILTERS * ACTIV_BITS)-1:0] conv2_out; wire conv2_valid; wire [((MFCC_FEATURES/2) * CONV2_NUM_FILTERS * ACTIV_BITS)-1:0] maxpool_out; wire maxpool_valid; wire [(FC1_OUTPUT_SIZE * ACTIV_BITS)-1:0] fc1_out; wire fc1_valid; wire [(FC2_OUTPUT_SIZE * ACTIV_BITS)-1:0] fc2_out; wire fc2_valid; wire [(FC2_OUTPUT_SIZE * ACTIV_BITS)-1:0] softmax_out; wire softmax_valid; // MFCC module instantiation mfcc_accel mfcc ( .clk(clk), .rst(rst_n), .audio_sample(audio_in), .sample_valid(audio_valid), .mfcc_feature(mfcc_out), // 640 bits .mfcc_valid(mfcc_valid) ); // Convolutional layer 1 conv2d #( .INPUT_WIDTH(1), // Single feature .INPUT_HEIGHT(MFCC_FEATURES), // 40 features .INPUT_CHANNELS(1), .KERNEL_SIZE(CONV1_KERNEL_SIZE), .NUM_FILTERS(CONV1_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv1 ( .clk(clk), .rst_n(rst_n), .data_in(mfcc_out), // 640 bits .data_valid(mfcc_valid), .data_out(conv1_out), .data_out_valid(conv1_valid), .weights_in(conv1_weights), .biases_in(conv1_biases), .load_weights(conv1_load_weights), .load_biases(conv1_load_biases) ); // Convolutional layer 2 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV1_NUM_FILTERS), .KERNEL_SIZE(CONV2_KERNEL_SIZE), .NUM_FILTERS(CONV2_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv2 ( .clk(clk), .rst_n(rst_n), .data_in(conv1_out), .data_valid(conv1_valid), .data_out(conv2_out), .data_out_valid(conv2_valid), .weights_in(conv2_weights), .biases_in(conv2_biases), .load_weights(conv2_load_weights), .load_biases(conv2_load_biases) ); // Maxpooling layer maxpool2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV2_NUM_FILTERS), .KERNEL_SIZE(2), .STRIDE(2), .ACTIV_BITS(ACTIV_BITS) ) maxpool ( .clk(clk), .rst_n(rst_n), .data_in(conv2_out), .data_valid(conv2_valid), .data_out(maxpool_out), .data_out_valid(maxpool_valid) ); // Fully connected layer 1 fully_connected #( .INPUT_SIZE(FC1_INPUT_SIZE), .OUTPUT_SIZE(FC1_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc1 ( .clk(clk), .rst_n(rst_n), .data_in(maxpool_out), .data_valid(maxpool_valid), .data_out(fc1_out), .data_out_valid(fc1_valid), .weights_in(fc1_weights), .biases_in(fc1_biases), .load_weights(fc1_load_weights), .load_biases(fc1_load_biases) ); // Fully connected layer 2 (output layer) fully_connected #( .INPUT_SIZE(FC2_INPUT_SIZE), .OUTPUT_SIZE(FC2_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc2 ( .clk(clk), .rst_n(rst_n), .data_in(fc1_out), .data_valid(fc1_valid), .data_out(fc2_out), .data_out_valid(fc2_valid), .weights_in(fc2_weights), .biases_in(fc2_biases), .load_weights(fc2_load_weights), .load_biases(fc2_load_biases) ); // Softmax layer softmax #( .INPUT_SIZE(NUM_KEYWORDS), .ACTIV_BITS(ACTIV_BITS) ) softmax ( .clk(clk), .rst_n(rst_n), .data_in(fc2_out), .data_valid(fc2_valid), .data_out(softmax_out), .data_out_valid(softmax_valid) ); // Output assignment always @(posedge clk or negedge rst_n) begin if (!rst_n) begin kws_result <= 'b0; kws_valid <= 1'b0; end else begin kws_result <= softmax_out[NUM_KEYWORDS*ACTIV_BITS-1:0]; kws_valid <= softmax_valid; end end endmodule

**Correct the Signal Widths in conv2d Module**

Ensure the **conv2d** module handles the input properly:

verilog

Copy code

module conv2d #( parameter INPUT_WIDTH = 40, // Number of MFCC features parameter INPUT_HEIGHT = 1, // Single feature height parameter INPUT_CHANNELS = 1, // Single channel input parameter KERNEL_SIZE = 3, parameter NUM_FILTERS = 8, parameter PADDING = 1, parameter ACTIV_BITS = 16 ) ( input wire clk, input wire rst_n, input wire [INPUT_WIDTH * INPUT_HEIGHT * INPUT_CHANNELS * ACTIV_BITS-1:0] data_in, input wire data_valid, output reg [INPUT_WIDTH *

INPUT_HEIGHT * NUM_FILTERS * ACTIV_BITS-1:0] data_out, output reg data_out_valid, input wire [NUM_FILTERS * INPUT_CHANNELS * KERNEL_SIZE * KERNEL_SIZE * ACTIV_BITS-1:0] weights_in, input wire [NUM_FILTERS * ACTIV_BITS-1:0] biases_in, input wire load_weights, input wire load_biases ); // Declare weights and biases reg [ACTIV_BITS-1:0] weights [0:NUM_FILTERS-1][0:INPUT_CHANNELS-1][0:KERNEL_SIZE-1][0:KERNEL_SIZE-1]; reg [ACTIV_BITS-1:0] biases [0:NUM_FILTERS-1]; // Declare internal signals reg [ACTIV_BITS-1:0] input_buffer [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:INPUT_CHANNELS-1]; reg [2*ACTIV_BITS-1:0] conv_result [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:NUM_FILTERS-1]; reg [ACTIV_BITS-1:0] relu_result [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:NUM_FILTERS-1]; // Load weights and biases integer i_load, j_load, k_load, l_load; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset weights and biases for (i_load = 0; i_load < NUM_FILTERS; i_load = i_load + 1) begin for (j_load = 0; j_load < INPUT_CHANNELS; j_load = j_load + 1) begin for (k_load = 0; k_load < KERNEL_SIZE; k_load = k_load + 1) begin for (l_load = 0; l_load < KERNEL_SIZE; l_load = l_load + 1) begin weights[i_load][j_load][k_load][l_load] <= 0; end end end biases[i_load] <= 0; end end else begin // Load weights when load_weights is asserted if (load_weights) begin for (i_load = 0; i_load < NUM_FILTERS; i_load = i_load + 1) begin for (j_load = 0; j_load < INPUT_CHANNELS; j_load = j_load + 1) begin for (k_load = 0; k_load < KERNEL_SIZE; k_load = k_load + 1) begin for (l_load = 0; l_load < KERNEL_SIZE; l_load = l_load + 1) begin weights[i_load][j_load][k_load][l_load] <= weights_in[(i_load*INPUT_CHANNELS*KERNEL_SIZE*KERNEL_SIZE + j_load*KERNEL_SIZE*KERNEL_SIZE + k_load*KERNEL_SIZE + l_load)*ACTIV_BITS +: ACTIV_BITS]; end end end end // Load biases when load_biases is asserted if (load_biases) begin for (i_load = 0; i_load < NUM_FILTERS; i_load = i_load + 1) begin biases[i_load] <= biases_in[i_load*ACTIV_BITS +: ACTIV_BITS]; end end end end // Convolution operation integer m_conv, n_conv, p_conv, q_conv, i_conv, j_conv, k_conv; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset internal signals and output for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < INPUT_CHANNELS; k_conv = k_conv + 1) begin input_buffer[i_conv][j_conv][k_conv] <= 0; end for (m_conv = 0; m_conv < NUM_FILTERS; m_conv = m_conv + 1) begin conv_result[i_conv][j_conv][m_conv] <= 0; relu_result[i_conv][j_conv][m_conv] <= 0; end end end data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Shift input data into buffer for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < INPUT_CHANNELS; k_conv = k_conv + 1) begin if (j_conv < INPUT_WIDTH - 1) begin input_buffer[i_conv][j_conv][k_conv] <= input_buffer[i_conv][j_conv+1][k_conv]; end else begin input_buffer[i_conv][j_conv][k_conv] <= data_in[i_conv*INPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS + j_conv*INPUT_CHANNELS*ACTIV_BITS + k_conv*ACTIV_BITS +: ACTIV_BITS]; end end end end // Perform convolution for (m_conv = 0; m_conv < INPUT_HEIGHT; m_conv = m_conv + 1) begin for (n_conv = 0; n_conv < INPUT_WIDTH; n_conv = n_conv + 1) begin for (p_conv = 0; p_conv < NUM_FILTERS; p_conv = p_conv + 1) begin conv_result[m_conv][n_conv][p_conv] = {{(2*ACTIV_BITS-ACTIV_BITS){1'b0}}, biases[p_conv]}; for (q_conv = 0; q_conv < INPUT_CHANNELS; q_conv = q_conv + 1) begin for (i_conv = 0; i_conv < KERNEL_SIZE; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < KERNEL_SIZE; j_conv = j_conv + 1) begin if (m_conv + i_conv - PADDING >= 0 && m_conv + i_conv - PADDING < INPUT_HEIGHT && n_conv + j_conv - PADDING >= 0 && n_conv + j_conv - PADDING < INPUT_WIDTH) begin conv_result[m_conv][n_conv][p_conv] = conv_result[m_conv][n_conv][p_conv] + weights[p_conv][q_conv][i_conv][j_conv] * input_buffer[m_conv + i_conv - PADDING][n_conv + j_conv - PADDING][q_conv]; end end end end end end end // Apply ReLU activation for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv <

NUM_FILTERS; k_conv = k_conv + 1) begin relu_result[i_conv][j_conv][k_conv] = (conv_result[i_conv][j_conv][k_conv][2*ACTIV_BITS-1] == 0) ? conv_result[i_conv][j_conv][k_conv][ACTIV_BITS-1:0] : 0; end end end // Assign output for (i_conv = 0; i_conv < INPUT_HEIGHT; i_conv = i_conv + 1) begin for (j_conv = 0; j_conv < INPUT_WIDTH; j_conv = j_conv + 1) begin for (k_conv = 0; k_conv < NUM_FILTERS; k_conv = k_conv + 1) begin data_out[i_conv*INPUT_WIDTH*NUM_FILTERS*ACTIV_BITS + j_conv*NUM_FILTERS*ACTIV_BITS + k_conv*ACTIV_BITS +: ACTIV_BITS] <= relu_result[i_conv][j_conv][k_conv]; end end end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule

**Update the maxpool2d Module**

verilog

Copy code

module maxpool2d #( parameter INPUT_WIDTH = 40, parameter INPUT_HEIGHT = 1, parameter INPUT_CHANNELS = 8, parameter KERNEL_SIZE = 2, parameter STRIDE = 2, parameter ACTIV_BITS = 16 ) ( input wire clk, input wire rst_n, input wire [INPUT_WIDTH * INPUT_HEIGHT * INPUT_CHANNELS * ACTIV_BITS-1:0] data_in, input wire data_valid, output reg [(INPUT_WIDTH/STRIDE) * INPUT_CHANNELS * ACTIV_BITS-1:0] data_out, output reg data_out_valid ); localparam OUTPUT_WIDTH = INPUT_WIDTH / STRIDE; localparam OUTPUT_HEIGHT = INPUT_HEIGHT / STRIDE; // Declare internal signals reg [ACTIV_BITS-1:0] input_buffer [0:INPUT_HEIGHT-1][0:INPUT_WIDTH-1][0:INPUT_CHANNELS-1]; reg [ACTIV_BITS-1:0] max_value [0:INPUT_CHANNELS-1]; // Max pooling operation integer i, j, k, m, n; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset internal signals and output for (i = 0; i < INPUT_HEIGHT; i = i + 1) begin for (j = 0; j < INPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin input_buffer[i][j][k] <= 0; end end end data_out <= 0; data_out_valid <= 0; end else if (data_valid) begin // Shift input data into buffer for (i = 0; i < INPUT_HEIGHT; i = i + 1) begin for (j = 0; j < INPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin if (j < INPUT_WIDTH - 1) begin input_buffer[i][j][k] <= input_buffer[i][j+1][k]; end else begin input_buffer[i][j][k] <= data_in[i*INPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS + j*INPUT_CHANNELS*ACTIV_BITS + k*ACTIV_BITS +: ACTIV_BITS]; end end end end // Perform max pooling for (i = 0; i < OUTPUT_HEIGHT; i = i + 1) begin for (j = 0; j < OUTPUT_WIDTH; j = j + 1) begin for (k = 0; k < INPUT_CHANNELS; k = k + 1) begin max_value[k] = input_buffer[i*STRIDE][j*STRIDE][k]; for (m = 0; m < KERNEL_SIZE; m = m + 1) begin for (n = 0; n < KERNEL_SIZE; n = n + 1) begin if (i*STRIDE + m < INPUT_HEIGHT && j*STRIDE + n < INPUT_WIDTH) begin max_value[k] = (input_buffer[i*STRIDE + m][j*STRIDE + n][k] > max_value[k]) ? input_buffer[i*STRIDE + m][j*STRIDE + n][k] : max_value[k]; end end end data_out[i*OUTPUT_WIDTH*INPUT_CHANNELS*ACTIV_BITS + j*INPUT_CHANNELS*ACTIV_BITS + k*ACTIV_BITS +: ACTIV_BITS] <= max_value[k]; end end end data_out_valid <= 1; end else begin data_out_valid <= 0; end end endmodule

**Conclusion**

By ensuring that each module correctly handles the 16-bit MFCC features and setting the signal widths accordingly, the design should now compile without width mismatch errors. The adjustments in the **mfcc_accel**, **conv2d**, and **maxpool2d** modules ensure proper signal width handling throughout the design.

Refactor the mfcc_accel. Create a reg [639:0] mfcc_feature (40 features * 16 bits = 640 bits). Create logic that will take the mfcc_feature signal and multiply it with features.

ChatGPT

Sure, let's refactor the **mfcc_accel** module to handle the 640-bit **mfcc_feature** signal and add logic to multiply it with features.

**Refactored mfcc_accel Module**

Here's how you can refactor the **mfcc_accel** module:

1.  Define a 640-bit register **mfcc_feature**.

2.  Add logic to handle the multiplication of the **mfcc_feature** signal with the features.

verilog

Copy code

module mfcc_accel ( input wire clk, input wire rst_n, input wire [15:0] audio_sample, input wire sample_valid, output reg [639:0] mfcc_feature, // 40 features * 16 bits = 640 bits output reg mfcc_valid ); // Declare internal registers and wires reg [15:0] feature_reg [0:39]; // Register array to hold 40 MFCC features integer i; always @(posedge clk or negedge rst_n) begin if (!rst_n) begin // Reset logic mfcc_valid <= 0; mfcc_feature <= 0; for (i = 0; i < 40; i = i + 1) begin feature_reg[i] <= 0; end end else if (sample_valid) begin // Assume feature extraction logic here and update the feature_reg for (i = 0; i < 40; i = i + 1) begin feature_reg[i] <= audio_sample; // Replace this with actual feature extraction logic end // Combine features into mfcc_feature mfcc_feature <= {feature_reg[39], feature_reg[38], feature_reg[37], feature_reg[36], feature_reg[35], feature_reg[34], feature_reg[33], feature_reg[32], feature_reg[31], feature_reg[30], feature_reg[29], feature_reg[28], feature_reg[27], feature_reg[26], feature_reg[25], feature_reg[24], feature_reg[23], feature_reg[22], feature_reg[21], feature_reg[20], feature_reg[19], feature_reg[18], feature_reg[17], feature_reg[16], feature_reg[15], feature_reg[14], feature_reg[13], feature_reg[12], feature_reg[11], feature_reg[10], feature_reg[9], feature_reg[8], feature_reg[7], feature_reg[6], feature_reg[5], feature_reg[4], feature_reg[3], feature_reg[2], feature_reg[1], feature_reg[0]}; // Signal that mfcc_feature is valid mfcc_valid <= 1; end else begin mfcc_valid <= 0; end end endmodule

**Explanation**

1.  **Register Definition**:

    *   A 640-bit register **mfcc_feature** is defined to hold the concatenated MFCC features.

    *   A register array **feature_reg** of 40 elements, each 16 bits wide, is used to store individual MFCC features.

2.  **Reset Logic**:

    *   On reset (**rst_n**), **mfcc_valid** and **mfcc_feature** are cleared.

    *   The **feature_reg** array is also cleared.

3.  **Feature Extraction and Combination**:

    *   When **sample_valid** is high, the module assumes some feature extraction logic (here, simplified as assigning **audio_sample** to all features for demonstration).

    *   The extracted features are combined into the **mfcc_feature** register.

- **mfcc_valid** is set to high to indicate that **mfcc_feature** is valid.

This code is a template and assumes a simplistic feature extraction logic for demonstration purposes. You should replace the feature extraction logic with actual MFCC computation.

Create a signal to load the weights. instantiate the rest of the modules "module mfcc_accel ( input wire clk, input wire rst, input wire [15:0] audio_sample, input wire sample_valid, output reg [15:0] mfcc_feature, // Adjusted width output reg mfcc_valid ); // Instantiate the Hanning window real module wire [15:0] hanning_real_out; wire hanning_real_valid; hanning_window_real hanning_real ( .clk(clk), .rst(rst), .sample_in(audio_sample), .sample_valid(sample_valid), .sample_out(hanning_real_out), .sample_out_valid(hanning_real_valid) ); // Instantiate the Hanning window imaginary module wire [15:0] hanning_imag_out; wire hanning_imag_valid; hanning_window_imag hanning_imag ( .clk(clk), .rst(rst), .sample_in(audio_sample), .sample_valid(sample_valid), .sample_out(hanning_imag_out), .sample_out_valid(hanning_imag_valid) ); // Instantiate the periodogram module wire [31:0] periodogram_out; wire periodogram_valid; periodogram_squared periodogram ( .clk(clk), .rst(rst), .sample_in_real(hanning_real_out), .sample_in_imag(hanning_imag_out), .sample_valid(hanning_real_valid), // Use the valid signal from the real Hanning window .periodogram_out(periodogram_out), .periodogram_valid(periodogram_valid) ); // Instantiate the pow_module wire [31:0] pow_out; wire pow_valid; pow_module pow ( .clk(clk), .rst(rst), .data_in(periodogram_out), .data_valid(periodogram_valid), .data_out(pow_out), .data_out_valid(pow_valid) ); // Instantiate the mel_filterbank module wire [31:0] mel_out; wire mel_valid; mel_filterbank mel ( .clk(clk), .rst(rst), .data_in(pow_out), .data_valid(pow_valid), .mel_out(mel_out), .mel_valid(mel_valid) ); // Instantiate the log_module wire [15:0] log_out; wire log_valid; log_module log ( .clk(clk), .rst(rst), .data_in(mel_out), .data_valid(mel_valid), .log_out(log_out), .log_valid(log_valid) ); // Instantiate the dct_module dct_module dct ( .clk(clk), .rst(rst), .data_in(log_out), .data_valid(log_valid), .dct_out(mfcc_feature), .dct_valid(mfcc_valid) ); endmodule"

ChatGPT

To refactor the **mfcc_accel** module and include instantiations of the submodules as described, we need to add the following:

- Define the module for handling the input and output ports.

- Instantiate the submodules (**hanning_window_real**, **hanning_window_imag**, **periodogram_squared**, **pow_module**, **mel_filterbank**, **log_module**, and **dct_module**).

- Define the signal to load the weights.

Here's how you can do this:

**Refactored mfcc_accel Module**

verilog

Copy code

module mfcc_accel ( input wire clk, input wire rst, input wire [15:0] audio_sample, input wire sample_valid, output reg [639:0] mfcc_feature, // 40 features * 16 bits = 640 bits output reg mfcc_valid ); // Signals for loading weights reg load_weights; wire [15:0] hanning_real_out; wire hanning_real_valid; // Instantiate the Hanning window real module hanning_window_real

hanning_real ( .clk(clk), .rst(rst), .sample_in(audio_sample), .sample_valid(sample_valid), .sample_out(hanning_real_out), .sample_out_valid(hanning_real_valid) ); // Instantiate the Hanning window imaginary module wire [15:0] hanning_imag_out; wire hanning_imag_valid; hanning_window_imag hanning_imag ( .clk(clk), .rst(rst), .sample_in(audio_sample), .sample_valid(sample_valid), .sample_out(hanning_imag_out), .sample_out_valid(hanning_imag_valid) ); // Instantiate the periodogram module wire [31:0] periodogram_out; wire periodogram_valid; periodogram_squared periodogram ( .clk(clk), .rst(rst), .sample_in_real(hanning_real_out), .sample_in_imag(hanning_imag_out), .sample_valid(hanning_real_valid), // Use the valid signal from the real Hanning window .periodogram_out(periodogram_out), .periodogram_valid(periodogram_valid) ); // Instantiate the pow_module wire [31:0] pow_out; wire pow_valid; pow_module pow ( .clk(clk), .rst(rst), .data_in(periodogram_out), .data_valid(periodogram_valid), .data_out(pow_out), .data_out_valid(pow_valid) ); // Instantiate the mel_filterbank module wire [31:0] mel_out; wire mel_valid; mel_filterbank mel ( .clk(clk), .rst(rst), .data_in(pow_out), .data_valid(pow_valid), .mel_out(mel_out), .mel_valid(mel_valid) ); // Instantiate the log_module wire [15:0] log_out; wire log_valid; log_module log ( .clk(clk), .rst(rst), .data_in(mel_out), .data_valid(mel_valid), .log_out(log_out), .log_valid(log_valid) ); // Instantiate the dct_module dct_module dct ( .clk(clk), .rst(rst), .data_in(log_out), .data_valid(log_valid), .dct_out(mfcc_feature), .dct_valid(mfcc_valid) ); endmodule

**Explanation**

1. **mfcc_accel Module**:

   - The module is designed to handle the MFCC feature extraction pipeline.

   - The **mfcc_feature** output is 640 bits wide, representing 40 MFCC features each 16 bits wide.

   - The **mfcc_valid** output signal indicates the validity of the **mfcc_feature**.

2. **Submodule Instantiations**:

   - **hanning_window_real** and **hanning_window_imag** modules process the audio samples to apply the Hanning window function.

   - **periodogram_squared** calculates the periodogram from the Hanning window outputs.

   - **pow_module** processes the periodogram to compute the power spectrum.

   - **mel_filterbank** applies the Mel filter bank to the power spectrum.

   - **log_module** calculates the logarithm of the Mel filter bank output.

   - **dct_module** performs the Discrete Cosine Transform (DCT) to obtain the MFCC features.

3. **Loading Weights Signal**:

   - The **load_weights** signal is defined but not used in this example. It can be integrated as needed for modules that require weight loading (e.g., neural network layers).

This setup ensures that the **mfcc_accel** module correctly processes the audio samples through the MFCC pipeline and outputs the MFCC features in the correct format.

refactor the top cnn_kws_accel based on these changes

ChatGPT

Let's refactor the **cnn_kws_accel** module to integrate the changes made to the **mfcc_accel** module. We need to ensure that the widths of the signals are consistent and correctly handle the outputs of the MFCC module through the CNN layers.

**Refactored cnn_kws_accel Module**

verilog

Copy code

```
module cnn_kws_accel #( parameter NUM_KEYWORDS = 10, parameter MFCC_FEATURES = 40,
parameter ACTIV_BITS = 16, // Each MFCC feature is 16 bits parameter FC1_INPUT_SIZE =
(MFCC_FEATURES/2) * (CONV2_NUM_FILTERS), parameter FC1_OUTPUT_SIZE = 64, parameter
FC2_INPUT_SIZE = 64, parameter FC2_OUTPUT_SIZE = NUM_KEYWORDS, parameter
CONV1_KERNEL_SIZE = 3, parameter CONV1_NUM_FILTERS = 8, parameter CONV2_KERNEL_SIZE = 3,
parameter CONV2_NUM_FILTERS = 8 )( input wire clk, input wire rst_n, input wire [15:0] audio_in,
input wire audio_valid, output reg [NUM_KEYWORDS-1:0] kws_result, output reg kws_valid, input
wire [7:0] frame_size, input wire [7:0] frame_overlap, input wire [7:0] num_mfcc_coeffs, input wire
[7:0] goertzel_coefs, output wire goertzel_coefs_start, output wire goertzel_coefs_valid, output wire
goertzel_coefs_done, input wire [CONV1_NUM_FILTERS * CONV1_KERNEL_SIZE *
CONV1_KERNEL_SIZE * ACTIV_BITS-1:0] conv1_weights, input wire [CONV1_NUM_FILTERS *
ACTIV_BITS-1:0] conv1_biases, input wire conv1_load_weights, input wire conv1_load_biases, input
wire [CONV2_NUM_FILTERS * CONV1_NUM_FILTERS * CONV2_KERNEL_SIZE * CONV2_KERNEL_SIZE
* ACTIV_BITS-1:0] conv2_weights, input wire [CONV2_NUM_FILTERS * ACTIV_BITS-1:0]
conv2_biases, input wire conv2_load_weights, input wire conv2_load_biases, input wire
[FC1_OUTPUT_SIZE * FC1_INPUT_SIZE * ACTIV_BITS-1:0] fc1_weights, input wire [FC1_OUTPUT_SIZE
* ACTIV_BITS-1:0] fc1_biases, input wire fc1_load_weights, input wire fc1_load_biases, input wire
[FC2_OUTPUT_SIZE * FC2_INPUT_SIZE * ACTIV_BITS-1:0] fc2_weights, input wire [FC2_OUTPUT_SIZE
* ACTIV_BITS-1:0] fc2_biases, input wire fc2_load_weights, input wire fc2_load_biases ); // MFCC
module signals wire [639:0] mfcc_out; // 40 features * 16 bits = 640 bits wire mfcc_valid; // CNN-
KWS layers wire [(MFCC_FEATURES * CONV1_NUM_FILTERS * ACTIV_BITS)-1:0] conv1_out; wire
conv1_valid; wire [(MFCC_FEATURES * CONV2_NUM_FILTERS * ACTIV_BITS)-1:0] conv2_out; wire
conv2_valid; wire [((MFCC_FEATURES/2) * CONV2_NUM_FILTERS * ACTIV_BITS)-1:0] maxpool_out;
wire maxpool_valid; wire [(FC1_OUTPUT_SIZE * ACTIV_BITS)-1:0] fc1_out; wire fc1_valid; wire
[(FC2_OUTPUT_SIZE * ACTIV_BITS)-1:0] fc2_out; wire fc2_valid; wire [(FC2_OUTPUT_SIZE *
ACTIV_BITS)-1:0] softmax_out; wire softmax_valid; // MFCC module instantiation mfcc_accel mfcc (
.clk(clk), .rst(rst_n), .audio_sample(audio_in), .sample_valid(audio_valid), .mfcc_feature(mfcc_out),
// 640 bits .mfcc_valid(mfcc_valid) ); // Convolutional layer 1 conv2d #(
.INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(1),
.KERNEL_SIZE(CONV1_KERNEL_SIZE), .NUM_FILTERS(CONV1_NUM_FILTERS), .PADDING(1),
.ACTIV_BITS(ACTIV_BITS) ) conv1 ( .clk(clk), .rst_n(rst_n), .data_in(mfcc_out), // 640 bits
.data_valid(mfcc_valid), .data_out(conv1_out), .data_out_valid(conv1_valid),
.weights_in(conv1_weights), .biases_in(conv1_biases), .load_weights(conv1_load_weights),
.load_biases(conv1_load_biases) ); // Convolutional layer 2 conv2d #(
```

.INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV1_NUM_FILTERS), .KERNEL_SIZE(CONV2_KERNEL_SIZE), .NUM_FILTERS(CONV2_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv2 ( .clk(clk), .rst_n(rst_n), .data_in(conv1_out), .data_valid(conv1_valid), .data_out(conv2_out), .data_out_valid(conv2_valid), .weights_in(conv2_weights), .biases_in(conv2_biases), .load_weights(conv2_load_weights), .load_biases(conv2_load_biases) ); // Maxpooling layer maxpool2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV2_NUM_FILTERS), .KERNEL_SIZE(2), .STRIDE(2), .ACTIV_BITS(ACTIV_BITS) ) maxpool ( .clk(clk), .rst_n(rst_n), .data_in(conv2_out), .data_valid(conv2_valid), .data_out(maxpool_out), .data_out_valid(maxpool_valid) ); // Fully connected layer 1 fully_connected #( .INPUT_SIZE(FC1_INPUT_SIZE), .OUTPUT_SIZE(FC1_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc1 ( .clk(clk), .rst_n(rst_n), .data_in(maxpool_out), .data_valid(maxpool_valid), .data_out(fc1_out), .data_out_valid(fc1_valid), .weights_in(fc1_weights), .biases_in(fc1_biases), .load_weights(fc1_load_weights), .load_biases(fc1_load_biases) ); // Fully connected layer 2 (output layer) fully_connected #( .INPUT_SIZE(FC2_INPUT_SIZE), .OUTPUT_SIZE(FC2_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc2 ( .clk(clk), .rst_n(rst_n), .data_in(fc1_out), .data_valid(fc1_valid), .data_out(fc2_out), .data_out_valid(fc2_valid), .weights_in(fc2_weights), .biases_in(fc2_biases), .load_weights(fc2_load_weights), .load_biases(fc2_load_biases) ); // Softmax layer softmax #( .INPUT_SIZE(NUM_KEYWORDS), .ACTIV_BITS(ACTIV_BITS) ) softmax ( .clk(clk), .rst_n(rst_n), .data_in(fc2_out), .data_valid(fc2_valid), .data_out(softmax_out), .data_out_valid(softmax_valid) ); // Output assignment always @(posedge clk or negedge rst_n) begin if (!rst_n) begin kws_result <= 'b0; kws_valid <= 1'b0; end else begin kws_result <= softmax_out[NUM_KEYWORDS*ACTIV_BITS-1:0]; kws_valid <= softmax_valid; end end endmodule

**Explanation**

1. **cnn_kws_accel Module**:

   - The module is designed to handle the entire KWS pipeline, from audio input to keyword spotting result.

   - The **mfcc_feature** output from the **mfcc_accel** module is 640 bits wide, representing 40 MFCC features each 16 bits wide.

2. **Submodule Instantiations**:

   - **mfcc_accel**: Handles the MFCC feature extraction and outputs the MFCC features and a valid signal.

   - **conv2d**: Handles the first and second convolutional layers. The input width and height are adjusted to match the MFCC feature dimensions.

   - **maxpool2d**: Handles max pooling after the second convolutional layer.

   - **fully_connected**: Handles the fully connected layers.

   - **softmax**: Handles the softmax layer to produce the final keyword spotting results.

3. **Weights and Biases**:

   - Inputs for weights and biases are included for each layer to allow for weight loading.

4. **Output Assignment**:

- The final keyword spotting results (**kws_result**) and validity (**kws_valid**) are assigned based on the softmax output.

This setup ensures that the **cnn_kws_accel** module correctly processes the audio samples through the MFCC extraction and CNN layers to output the keyword spotting results.

The mfcc_accel module "module mfcc_accel ( input wire clk, input wire rst, input wire [15:0] audio_sample, input wire sample_valid, output reg [639:0] mfcc_feature, // 40 features * 16 bits = 640 bits output reg mfcc_valid ); // Signals for loading weights reg load_weights; wire [15:0] hanning_real_out; wire hanning_real_valid; // Instantiate the Hanning window real module hanning_window_real hanning_real ( .clk(clk), .rst(rst), .sample_in(audio_sample), .sample_valid(sample_valid), .sample_out(hanning_real_out), .sample_out_valid(hanning_real_valid) ); // Instantiate the Hanning window imaginary module wire [15:0] hanning_imag_out; wire hanning_imag_valid; hanning_window_imag hanning_imag ( .clk(clk), .rst(rst), .sample_in(audio_sample), .sample_valid(sample_valid), .sample_out(hanning_imag_out), .sample_out_valid(hanning_imag_valid) ); // Instantiate the periodogram module wire [31:0] periodogram_out; wire periodogram_valid; periodogram_squared periodogram ( .clk(clk), .rst(rst), .sample_in_real(hanning_real_out), .sample_in_imag(hanning_imag_out), .sample_valid(hanning_real_valid), // Use the valid signal from the real Hanning window .periodogram_out(periodogram_out), .periodogram_valid(periodogram_valid) ); // Instantiate the pow_module wire [31:0] pow_out; wire pow_valid; pow_module pow ( .clk(clk), .rst(rst), .data_in(periodogram_out), .data_valid(periodogram_valid), .data_out(pow_out), .data_out_valid(pow_valid) ); // Instantiate the mel_filterbank module wire [31:0] mel_out; wire mel_valid; mel_filterbank mel ( .clk(clk), .rst(rst), .data_in(pow_out), .data_valid(pow_valid), .mel_out(mel_out), .mel_valid(mel_valid) ); // Instantiate the log_module wire [15:0] log_out; wire log_valid; log_module log ( .clk(clk), .rst(rst), .data_in(mel_out), .data_valid(mel_valid), .log_out(log_out), .log_valid(log_valid) ); // Instantiate the dct_module dct_module dct ( .clk(clk), .rst(rst), .data_in(log_out), .data_valid(log_valid), .dct_out(mfcc_feature), .dct_valid(mfcc_valid) ); endmodule " has the following submodules: "module i2c_slave ( // I2C interface signals input wire scl_i, input wire sda_i, output reg sda_o, output reg sda_t, // Output audio signal output reg [15:0] audio_data_o, output reg audio_valid_o ); // I2C slave address parameter I2C_SLAVE_ADDRESS = 7'h45; // I2C slave state machine parameter STATE_IDLE = 3'b000; parameter STATE_ADDR = 3'b001; parameter STATE_ACK_ADDR = 3'b010; parameter STATE_READ = 3'b011; parameter STATE_ACK_DATA = 3'b100; parameter STATE_STOP = 3'b101; // Internal registers and signals reg [2:0] i2c_state; reg [6:0] i2c_slave_addr; reg i2c_rw; reg [3:0] i2c_bit_count; reg [15:0] i2c_data_reg; reg i2c_data_valid; wire i2c_start; wire i2c_stop; wire i2c_ack; // I2C start and stop condition detection assign i2c_start = !sda_i && scl_i; assign i2c_stop = sda_i && scl_i; // I2C acknowledge detection assign i2c_ack = !sda_i; // I2C slave logic always @(posedge scl_i or negedge scl_i) begin case (i2c_state) STATE_IDLE: begin if (i2c_start) begin i2c_bit_count <= 7; i2c_state <= STATE_ADDR; audio_valid_o <= 0; end end STATE_ADDR: begin if (scl_i) begin i2c_slave_addr <= {i2c_slave_addr[5:0], sda_i}; i2c_bit_count <= i2c_bit_count - 1; if (i2c_bit_count == 0) begin i2c_rw <= sda_i; i2c_state <= STATE_ACK_ADDR; end end end STATE_ACK_ADDR: begin if (!scl_i) begin if (i2c_slave_addr == I2C_SLAVE_ADDRESS[6:0]) begin sda_o <= 0; // ACK if (i2c_rw) begin i2c_bit_count <= 15; i2c_state <= STATE_READ; end else begin i2c_state <= STATE_IDLE; end end else begin sda_o <= 1; // NACK i2c_state <= STATE_IDLE; end end else begin sda_o <= 1; // Release SDA end end STATE_READ: begin if (scl_i) begin i2c_data_reg <= {i2c_data_reg[14:0], sda_i}; i2c_bit_count <= i2c_bit_count - 1; if (i2c_bit_count == 0) begin i2c_state <= STATE_ACK_DATA; end end end STATE_ACK_DATA: begin if (!scl_i) begin sda_o <= 0; // ACK i2c_data_valid <= 1; audio_data_o <= i2c_data_reg; audio_valid_o <= 1; end else begin sda_o <= 1; // Release SDA if

(i2c_stop) begin i2c_state <= STATE_STOP; end else begin i2c_bit_count <= 15; i2c_state <= STATE_READ; end end end STATE_STOP: begin i2c_state <= STATE_IDLE; i2c_data_valid <= 0; audio_valid_o <= 0; end default: begin i2c_state <= STATE_IDLE; audio_valid_o <= 0; end endcase end // I2C data and clock control always @(*) begin case (i2c_state) STATE_ACK_ADDR, STATE_ACK_DATA: begin sda_t = 0; // Drive SDA low for ACK end default: begin sda_t = 1; // Release SDA end endcase end endmodule ", "// Hanning Window Imaginary Module module hanning_window_imag ( input wire clk, input wire rst, input wire [15:0] sample_in, input wire sample_valid, output reg [15:0] sample_out, output reg sample_out_valid ); // Parameters localparam N = 256; // Frame size localparam Q = 15; // Fixed-point precision localparam NF = 512; // Power-of-two size for zero-padding // Registers reg [15:0] sample_buffer [0:N-1]; reg [$clog2(NF)-1:0] sample_count; reg [$clog2(N)-1:0] coeff_count; reg [15:0] coeff; // Fixed-point constants localparam [15:0] CONST_05 = 16'h4000; // 0.5 in Q15 localparam [15:0] CONST_2PI = 16'h6487; // 2π in Q15 // Hanning window coefficients (pre-computed) reg [15:0] hanning_coeff [0:N-1]; initial begin hanning_coeff[0] = 16'h0000; hanning_coeff[1] = 16'h0005; hanning_coeff[2] = 16'h0014; hanning_coeff[3] = 16'h002D; hanning_coeff[4] = 16'h0050; hanning_coeff[5] = 16'h007C; hanning_coeff[6] = 16'h00B3; hanning_coeff[7] = 16'h00F3; hanning_coeff[8] = 16'h013D; hanning_coeff[9] = 16'h0191; hanning_coeff[10] = 16'h01EF; hanning_coeff[11] = 16'h0256; hanning_coeff[12] = 16'h02C7; hanning_coeff[13] = 16'h0341; hanning_coeff[14] = 16'h03C5; hanning_coeff[15] = 16'h0452; hanning_coeff[16] = 16'h04E9; hanning_coeff[17] = 16'h0588; hanning_coeff[18] = 16'h0631; hanning_coeff[19] = 16'h06E3; hanning_coeff[20] = 16'h079D; hanning_coeff[21] = 16'h0861; hanning_coeff[22] = 16'h092D; hanning_coeff[23] = 16'h0A01; hanning_coeff[24] = 16'h0ADE; hanning_coeff[25] = 16'h0BC3; hanning_coeff[26] = 16'h0CB1; hanning_coeff[27] = 16'h0DA6; hanning_coeff[28] = 16'h0EA3; hanning_coeff[29] = 16'h0FA8; hanning_coeff[30] = 16'h10B4; hanning_coeff[31] = 16'h11C8; hanning_coeff[32] = 16'h12E2; hanning_coeff[33] = 16'h1404; hanning_coeff[34] = 16'h152D; hanning_coeff[35] = 16'h165C; hanning_coeff[36] = 16'h1792; hanning_coeff[37] = 16'h18CE; hanning_coeff[38] = 16'h1A10; hanning_coeff[39] = 16'h1B58; hanning_coeff[40] = 16'h1CA6; hanning_coeff[41] = 16'h1DF9; hanning_coeff[42] = 16'h1F52; hanning_coeff[43] = 16'h20AF; hanning_coeff[44] = 16'h2212; hanning_coeff[45] = 16'h2379; hanning_coeff[46] = 16'h24E5; hanning_coeff[47] = 16'h2654; hanning_coeff[48] = 16'h27C8; hanning_coeff[49] = 16'h2940; hanning_coeff[50] = 16'h2ABB; hanning_coeff[51] = 16'h2C39; hanning_coeff[52] = 16'h2DBB; hanning_coeff[53] = 16'h2F3F; hanning_coeff[54] = 16'h30C6; hanning_coeff[55] = 16'h324F; hanning_coeff[56] = 16'h33DA; hanning_coeff[57] = 16'h3568; hanning_coeff[58] = 16'h36F6; hanning_coeff[59] = 16'h3887; hanning_coeff[60] = 16'h3A18; hanning_coeff[61] = 16'h3BAB; hanning_coeff[62] = 16'h3D3E; hanning_coeff[63] = 16'h3ED1; hanning_coeff[64] = 16'h4065; hanning_coeff[65] = 16'h41F9; hanning_coeff[66] = 16'h438C; hanning_coeff[67] = 16'h451F; hanning_coeff[68] = 16'h46B1; hanning_coeff[69] = 16'h4842; hanning_coeff[70] = 16'h49D1; hanning_coeff[71] = 16'h4B5F; hanning_coeff[72] = 16'h4CEC; hanning_coeff[73] = 16'h4E76; hanning_coeff[74] = 16'h4FFE; hanning_coeff[75] = 16'h5184; hanning_coeff[76] = 16'h5307; hanning_coeff[77] = 16'h5487; hanning_coeff[78] = 16'h5603; hanning_coeff[79] = 16'h577D; hanning_coeff[80] = 16'h58F2; hanning_coeff[81] = 16'h5A64; hanning_coeff[82] = 16'h5BD2; hanning_coeff[83] = 16'h5D3B; hanning_coeff[84] = 16'h5EA0; hanning_coeff[85] = 16'h6000; hanning_coeff[86] = 16'h615B; hanning_coeff[87] = 16'h62B1; hanning_coeff[88] = 16'h6402; hanning_coeff[89] = 16'h654C; hanning_coeff[90] = 16'h6692; hanning_coeff[91] = 16'h67D1; hanning_coeff[92] = 16'h690A; hanning_coeff[93] = 16'h6A3C; hanning_coeff[94] = 16'h6B68; hanning_coeff[95] = 16'h6C8D; hanning_coeff[96] = 16'h6DAC; hanning_coeff[97] = 16'h6EC3; hanning_coeff[98] = 16'h6FD3; hanning_coeff[99] = 16'h70DC; hanning_coeff[100] = 16'h71DC; hanning_coeff[101] = 16'h72D6; hanning_coeff[102] = 16'h73C7; hanning_coeff[103] = 16'h74B0; hanning_coeff[104] = 16'h7591; hanning_coeff[105] = 16'h766A;

```
hanning_coeff[106] = 16'h773A; hanning_coeff[107] = 16'h7802; hanning_coeff[108] = 16'h78C1;
hanning_coeff[109] = 16'h7977; hanning_coeff[110] = 16'h7A24; hanning_coeff[111] = 16'h7AC8;
hanning_coeff[112] = 16'h7B64; hanning_coeff[113] = 16'h7BF5; hanning_coeff[114] = 16'h7C7E;
hanning_coeff[115] = 16'h7CFD; hanning_coeff[116] = 16'h7D73; hanning_coeff[117] = 16'h7DDF;
hanning_coeff[118] = 16'h7E41; hanning_coeff[119] = 16'h7E9A; hanning_coeff[120] = 16'h7EE9;
hanning_coeff[121] = 16'h7F2E; hanning_coeff[122] = 16'h7F6A; hanning_coeff[123] = 16'h7F9B;
hanning_coeff[124] = 16'h7FC3; hanning_coeff[125] = 16'h7FE1; hanning_coeff[126] = 16'h7FF5;
hanning_coeff[127] = 16'h7FFF; hanning_coeff[128] = 16'h7FFF; hanning_coeff[129] = 16'h7FF5;
hanning_coeff[130] = 16'h7FE1; hanning_coeff[131] = 16'h7FC3; hanning_coeff[132] = 16'h7F9B;
hanning_coeff[133] = 16'h7F6A; hanning_coeff[134] = 16'h7F2E; hanning_coeff[135] = 16'h7EE9;
hanning_coeff[136] = 16'h7E9A; hanning_coeff[137] = 16'h7E41; hanning_coeff[138] = 16'h7DDF;
hanning_coeff[139] = 16'h7D73; hanning_coeff[140] = 16'h7CFD; hanning_coeff[141] = 16'h7C7E;
hanning_coeff[142] = 16'h7BF5; hanning_coeff[143] = 16'h7B64; hanning_coeff[144] = 16'h7AC8;
hanning_coeff[145] = 16'h7A24; hanning_coeff[146] = 16'h7977; hanning_coeff[147] = 16'h78C1;
hanning_coeff[148] = 16'h7802; hanning_coeff[149] = 16'h773A; hanning_coeff[150] = 16'h766A;
hanning_coeff[151] = 16'h7591; hanning_coeff[152] = 16'h74B0; hanning_coeff[153] = 16'h73C7;
hanning_coeff[154] = 16'h72D6; hanning_coeff[155] = 16'h71DC; hanning_coeff[156] = 16'h70DC;
hanning_coeff[157] = 16'h6FD3; hanning_coeff[158] = 16'h6EC3; hanning_coeff[159] = 16'h6DAC;
hanning_coeff[160] = 16'h6C8D; hanning_coeff[161] = 16'h6B68; hanning_coeff[162] = 16'h6A3C;
hanning_coeff[163] = 16'h690A; hanning_coeff[164] = 16'h67D1; hanning_coeff[165] = 16'h6692;
hanning_coeff[166] = 16'h654C; hanning_coeff[167] = 16'h6402; hanning_coeff[168] = 16'h62B1;
hanning_coeff[169] = 16'h615B; hanning_coeff[170] = 16'h6000; hanning_coeff[171] = 16'h5EA0;
hanning_coeff[172] = 16'h5D3B; hanning_coeff[173] = 16'h5BD2; hanning_coeff[174] = 16'h5A64;
hanning_coeff[175] = 16'h58F2; hanning_coeff[176] = 16'h577D; hanning_coeff[177] = 16'h5603;
hanning_coeff[178] = 16'h5487; hanning_coeff[179] = 16'h5307; hanning_coeff[180] = 16'h5184;
hanning_coeff[181] = 16'h4FFE; hanning_coeff[182] = 16'h4E76; hanning_coeff[183] = 16'h4CEC;
hanning_coeff[184] = 16'h4B5F; hanning_coeff[185] = 16'h49D1; hanning_coeff[186] = 16'h4842;
hanning_coeff[187] = 16'h46B1; hanning_coeff[188] = 16'h451F; hanning_coeff[189] = 16'h438C;
hanning_coeff[190] = 16'h41F9; hanning_coeff[191] = 16'h4065; hanning_coeff[192] = 16'h3ED1;
hanning_coeff[193] = 16'h3D3E; hanning_coeff[194] = 16'h3BAB; hanning_coeff[195] = 16'h3A18;
hanning_coeff[196] = 16'h3887; hanning_coeff[197] = 16'h36F6; hanning_coeff[198] = 16'h3568;
hanning_coeff[199] = 16'h33DA; hanning_coeff[200] = 16'h324F; hanning_coeff[201] = 16'h30C6;
hanning_coeff[202] = 16'h2F3F; hanning_coeff[203] = 16'h2DBB; hanning_coeff[204] = 16'h2C39;
hanning_coeff[205] = 16'h2ABB; hanning_coeff[206] = 16'h2940; hanning_coeff[207] = 16'h27C8;
hanning_coeff[208] = 16'h2654; hanning_coeff[209] = 16'h24E5; hanning_coeff[210] = 16'h2379;
hanning_coeff[211] = 16'h2212; hanning_coeff[212] = 16'h20AF; hanning_coeff[213] = 16'h1F52;
hanning_coeff[214] = 16'h1DF9; hanning_coeff[215] = 16'h1CA6; hanning_coeff[216] = 16'h1B58;
hanning_coeff[217] = 16'h1A10; hanning_coeff[218] = 16'h18CE; hanning_coeff[219] = 16'h1792;
hanning_coeff[220] = 16'h165C; hanning_coeff[221] = 16'h152D; hanning_coeff[222] = 16'h1404;
hanning_coeff[223] = 16'h12E2; hanning_coeff[224] = 16'h11C8; hanning_coeff[225] = 16'h10B4;
hanning_coeff[226] = 16'h0FA8; hanning_coeff[227] = 16'h0EA3; hanning_coeff[228] = 16'h0DA6;
hanning_coeff[229] = 16'h0CB1; hanning_coeff[230] = 16'h0BC3; hanning_coeff[231] = 16'h0ADE;
hanning_coeff[232] = 16'h0A01; hanning_coeff[233] = 16'h092D; hanning_coeff[234] = 16'h0861;
hanning_coeff[235] = 16'h079D; hanning_coeff[236] = 16'h06E3; hanning_coeff[237] = 16'h0631;
hanning_coeff[238] = 16'h0588; hanning_coeff[239] = 16'h04E9; hanning_coeff[240] = 16'h0452;
hanning_coeff[241] = 16'h03C5; hanning_coeff[242] = 16'h0341; hanning_coeff[243] = 16'h02C7;
hanning_coeff[244] = 16'h0256; hanning_coeff[245] = 16'h01EF; hanning_coeff[246] = 16'h0191;
hanning_coeff[247] = 16'h013D; hanning_coeff[248] = 16'h00F3; hanning_coeff[249] = 16'h00B3;
```

hanning_coeff[250] = 16'h007C; hanning_coeff[251] = 16'h0050; hanning_coeff[252] = 16'h002D; hanning_coeff[253] = 16'h0014; hanning_coeff[254] = 16'h0005; hanning_coeff[255] = 16'h0000; end // Input sample buffer always @(posedge clk) begin if (rst) begin sample_count <= 0; end else if (sample_valid) begin sample_buffer[sample_count] <= sample_in; sample_count <= (sample_count == N-1) ? 0 : sample_count + 1; end end // Hanning window processing always @(posedge clk) begin if (rst) begin coeff_count <= 0; sample_out <= 0; sample_out_valid <= 0; end else begin if (sample_count == N-1 && sample_valid) begin coeff_count <= 0; sample_out_valid <= 1; end else if (coeff_count < N) begin coeff <= hanning_coeff[coeff_count]; sample_out <= (sample_buffer[coeff_count] * coeff) >>> Q; coeff_count <= coeff_count + 1; end else if (coeff_count < NF) begin sample_out <= 0; // Zero-padding coeff_count <= coeff_count + 1; end else begin sample_out_valid <= 0; end end end endmodule ","// Hanning Window Real Module module hanning_window_real ( input wire clk, input wire rst, input wire [15:0] sample_in, input wire sample_valid, output reg [15:0] sample_out, output reg sample_out_valid ); // Parameters localparam N = 256; // Frame size localparam Q = 15; // Fixed-point precision localparam NF = 512; // Power-of-two size for zero-padding // Registers reg [15:0] sample_buffer [0:N-1]; reg [$clog2(NF)-1:0] sample_count; reg [$clog2(N)-1:0] coeff_count; reg [15:0] coeff; // Fixed-point constants localparam [15:0] CONST_05 = 16'h4000; // 0.5 in Q15 localparam [15:0] CONST_2PI = 16'h6487; // 2π in Q15 // Hanning window coefficients (pre-computed) reg [15:0] hanning_coeff [0:N-1]; initial begin
hanning_coeff[0] = 16'h0000; hanning_coeff[1] = 16'h0005; hanning_coeff[2] = 16'h0014;
hanning_coeff[3] = 16'h002D; hanning_coeff[4] = 16'h0050; hanning_coeff[5] = 16'h007C;
hanning_coeff[6] = 16'h00B3; hanning_coeff[7] = 16'h00F3; hanning_coeff[8] = 16'h013D;
hanning_coeff[9] = 16'h0191; hanning_coeff[10] = 16'h01EF; hanning_coeff[11] = 16'h0256;
hanning_coeff[12] = 16'h02C7; hanning_coeff[13] = 16'h0341; hanning_coeff[14] = 16'h03C5;
hanning_coeff[15] = 16'h0452; hanning_coeff[16] = 16'h04E9; hanning_coeff[17] = 16'h0588;
hanning_coeff[18] = 16'h0631; hanning_coeff[19] = 16'h06E3; hanning_coeff[20] = 16'h079D;
hanning_coeff[21] = 16'h0861; hanning_coeff[22] = 16'h092D; hanning_coeff[23] = 16'h0A01;
hanning_coeff[24] = 16'h0ADE; hanning_coeff[25] = 16'h0BC3; hanning_coeff[26] = 16'h0CB1;
hanning_coeff[27] = 16'h0DA6; hanning_coeff[28] = 16'h0EA3; hanning_coeff[29] = 16'h0FA8;
hanning_coeff[30] = 16'h10B4; hanning_coeff[31] = 16'h11C8; hanning_coeff[32] = 16'h12E2;
hanning_coeff[33] = 16'h1404; hanning_coeff[34] = 16'h152D; hanning_coeff[35] = 16'h165C;
hanning_coeff[36] = 16'h1792; hanning_coeff[37] = 16'h18CE; hanning_coeff[38] = 16'h1A10;
hanning_coeff[39] = 16'h1B58; hanning_coeff[40] = 16'h1CA6; hanning_coeff[41] = 16'h1DF9;
hanning_coeff[42] = 16'h1F52; hanning_coeff[43] = 16'h20AF; hanning_coeff[44] = 16'h2212;
hanning_coeff[45] = 16'h2379; hanning_coeff[46] = 16'h24E5; hanning_coeff[47] = 16'h2654;
hanning_coeff[48] = 16'h27C8; hanning_coeff[49] = 16'h2940; hanning_coeff[50] = 16'h2ABB;
hanning_coeff[51] = 16'h2C39; hanning_coeff[52] = 16'h2DBB; hanning_coeff[53] = 16'h2F3F;
hanning_coeff[54] = 16'h30C6; hanning_coeff[55] = 16'h324F; hanning_coeff[56] = 16'h33DA;
hanning_coeff[57] = 16'h3568; hanning_coeff[58] = 16'h36F6; hanning_coeff[59] = 16'h3887;
hanning_coeff[60] = 16'h3A18; hanning_coeff[61] = 16'h3BAB; hanning_coeff[62] = 16'h3D3E;
hanning_coeff[63] = 16'h3ED1; hanning_coeff[64] = 16'h4065; hanning_coeff[65] = 16'h41F9;
hanning_coeff[66] = 16'h438C; hanning_coeff[67] = 16'h451F; hanning_coeff[68] = 16'h46B1;
hanning_coeff[69] = 16'h4842; hanning_coeff[70] = 16'h49D1; hanning_coeff[71] = 16'h4B5F;
hanning_coeff[72] = 16'h4CEC; hanning_coeff[73] = 16'h4E76; hanning_coeff[74] = 16'h4FFE;
hanning_coeff[75] = 16'h5184; hanning_coeff[76] = 16'h5307; hanning_coeff[77] = 16'h5487;
hanning_coeff[78] = 16'h5603; hanning_coeff[79] = 16'h577D; hanning_coeff[80] = 16'h58F2;
hanning_coeff[81] = 16'h5A64; hanning_coeff[82] = 16'h5BD2; hanning_coeff[83] = 16'h5D3B;
hanning_coeff[84] = 16'h5EA0; hanning_coeff[85] = 16'h6000; hanning_coeff[86] = 16'h615B;
hanning_coeff[87] = 16'h62B1; hanning_coeff[88] = 16'h6402; hanning_coeff[89] = 16'h654C;

hanning_coeff[90] = 16'h6692; hanning_coeff[91] = 16'h67D1; hanning_coeff[92] = 16'h690A;
hanning_coeff[93] = 16'h6A3C; hanning_coeff[94] = 16'h6B68; hanning_coeff[95] = 16'h6C8D;
hanning_coeff[96] = 16'h6DAC; hanning_coeff[97] = 16'h6EC3; hanning_coeff[98] = 16'h6FD3;
hanning_coeff[99] = 16'h70DC; hanning_coeff[100] = 16'h71DC; hanning_coeff[101] = 16'h72D6;
hanning_coeff[102] = 16'h73C7; hanning_coeff[103] = 16'h74B0; hanning_coeff[104] = 16'h7591;
hanning_coeff[105] = 16'h766A; hanning_coeff[106] = 16'h773A; hanning_coeff[107] = 16'h7802;
hanning_coeff[108] = 16'h78C1; hanning_coeff[109] = 16'h7977; hanning_coeff[110] = 16'h7A24;
hanning_coeff[111] = 16'h7AC8; hanning_coeff[112] = 16'h7B64; hanning_coeff[113] = 16'h7BF5;
hanning_coeff[114] = 16'h7C7E; hanning_coeff[115] = 16'h7CFD; hanning_coeff[116] = 16'h7D73;
hanning_coeff[117] = 16'h7DDF; hanning_coeff[118] = 16'h7E41; hanning_coeff[119] = 16'h7E9A;
hanning_coeff[120] = 16'h7EE9; hanning_coeff[121] = 16'h7F2E; hanning_coeff[122] = 16'h7F6A;
hanning_coeff[123] = 16'h7F9B; hanning_coeff[124] = 16'h7FC3; hanning_coeff[125] = 16'h7FE1;
hanning_coeff[126] = 16'h7FF5; hanning_coeff[127] = 16'h7FFF; hanning_coeff[128] = 16'h7FFF;
hanning_coeff[129] = 16'h7FF5; hanning_coeff[130] = 16'h7FE1; hanning_coeff[131] = 16'h7FC3;
hanning_coeff[132] = 16'h7F9B; hanning_coeff[133] = 16'h7F6A; hanning_coeff[134] = 16'h7F2E;
hanning_coeff[135] = 16'h7EE9; hanning_coeff[136] = 16'h7E9A; hanning_coeff[137] = 16'h7E41;
hanning_coeff[138] = 16'h7DDF; hanning_coeff[139] = 16'h7D73; hanning_coeff[140] = 16'h7CFD;
hanning_coeff[141] = 16'h7C7E; hanning_coeff[142] = 16'h7BF5; hanning_coeff[143] = 16'h7B64;
hanning_coeff[144] = 16'h7AC8; hanning_coeff[145] = 16'h7A24; hanning_coeff[146] = 16'h7977;
hanning_coeff[147] = 16'h78C1; hanning_coeff[148] = 16'h7802; hanning_coeff[149] = 16'h773A;
hanning_coeff[150] = 16'h766A; hanning_coeff[151] = 16'h7591; hanning_coeff[152] = 16'h74B0;
hanning_coeff[153] = 16'h73C7; hanning_coeff[154] = 16'h72D6; hanning_coeff[155] = 16'h71DC;
hanning_coeff[156] = 16'h70DC; hanning_coeff[157] = 16'h6FD3; hanning_coeff[158] = 16'h6EC3;
hanning_coeff[159] = 16'h6DAC; hanning_coeff[160] = 16'h6C8D; hanning_coeff[161] = 16'h6B68;
hanning_coeff[162] = 16'h6A3C; hanning_coeff[163] = 16'h690A; hanning_coeff[164] = 16'h67D1;
hanning_coeff[165] = 16'h6692; hanning_coeff[166] = 16'h654C; hanning_coeff[167] = 16'h6402;
hanning_coeff[168] = 16'h62B1; hanning_coeff[169] = 16'h615B; hanning_coeff[170] = 16'h6000;
hanning_coeff[171] = 16'h5EA0; hanning_coeff[172] = 16'h5D3B; hanning_coeff[173] = 16'h5BD2;
hanning_coeff[174] = 16'h5A64; hanning_coeff[175] = 16'h58F2; hanning_coeff[176] = 16'h577D;
hanning_coeff[177] = 16'h5603; hanning_coeff[178] = 16'h5487; hanning_coeff[179] = 16'h5307;
hanning_coeff[180] = 16'h5184; hanning_coeff[181] = 16'h4FFE; hanning_coeff[182] = 16'h4E76;
hanning_coeff[183] = 16'h4CEC; hanning_coeff[184] = 16'h4B5F; hanning_coeff[185] = 16'h49D1;
hanning_coeff[186] = 16'h4842; hanning_coeff[187] = 16'h46B1; hanning_coeff[188] = 16'h451F;
hanning_coeff[189] = 16'h438C; hanning_coeff[190] = 16'h41F9; hanning_coeff[191] = 16'h4065;
hanning_coeff[192] = 16'h3ED1; hanning_coeff[193] = 16'h3D3E; hanning_coeff[194] = 16'h3BAB;
hanning_coeff[195] = 16'h3A18; hanning_coeff[196] = 16'h3887; hanning_coeff[197] = 16'h36F6;
hanning_coeff[198] = 16'h3568; hanning_coeff[199] = 16'h33DA; hanning_coeff[200] = 16'h324F;
hanning_coeff[201] = 16'h30C6; hanning_coeff[202] = 16'h2F3F; hanning_coeff[203] = 16'h2DBB;
hanning_coeff[204] = 16'h2C39; hanning_coeff[205] = 16'h2ABB; hanning_coeff[206] = 16'h2940;
hanning_coeff[207] = 16'h27C8; hanning_coeff[208] = 16'h2654; hanning_coeff[209] = 16'h24E5;
hanning_coeff[210] = 16'h2379; hanning_coeff[211] = 16'h2212; hanning_coeff[212] = 16'h20AF;
hanning_coeff[213] = 16'h1F52; hanning_coeff[214] = 16'h1DF9; hanning_coeff[215] = 16'h1CA6;
hanning_coeff[216] = 16'h1B58; hanning_coeff[217] = 16'h1A10; hanning_coeff[218] = 16'h18CE;
hanning_coeff[219] = 16'h1792; hanning_coeff[220] = 16'h165C; hanning_coeff[221] = 16'h152D;
hanning_coeff[222] = 16'h1404; hanning_coeff[223] = 16'h12E2; hanning_coeff[224] = 16'h11C8;
hanning_coeff[225] = 16'h10B4; hanning_coeff[226] = 16'h0FA8; hanning_coeff[227] = 16'h0EA3;
hanning_coeff[228] = 16'h0DA6; hanning_coeff[229] = 16'h0CB1; hanning_coeff[230] = 16'h0BC3;
hanning_coeff[231] = 16'h0ADE; hanning_coeff[232] = 16'h0A01; hanning_coeff[233] = 16'h092D;

hanning_coeff[234] = 16'h0861; hanning_coeff[235] = 16'h079D; hanning_coeff[236] = 16'h06E3; hanning_coeff[237] = 16'h0631; hanning_coeff[238] = 16'h0588; hanning_coeff[239] = 16'h04E9; hanning_coeff[240] = 16'h0452; hanning_coeff[241] = 16'h03C5; hanning_coeff[242] = 16'h0341; hanning_coeff[243] = 16'h02C7; hanning_coeff[244] = 16'h0256; hanning_coeff[245] = 16'h01EF; hanning_coeff[246] = 16'h0191; hanning_coeff[247] = 16'h013D; hanning_coeff[248] = 16'h00F3; hanning_coeff[249] = 16'h00B3; hanning_coeff[250] = 16'h007C; hanning_coeff[251] = 16'h0050; hanning_coeff[252] = 16'h002D; hanning_coeff[253] = 16'h0014; hanning_coeff[254] = 16'h0005; hanning_coeff[255] = 16'h0000; end // Input sample buffer always @(posedge clk) begin if (rst) begin sample_count <= 0; end else if (sample_valid) begin sample_buffer[sample_count] <= sample_in; sample_count <= (sample_count == N-1) ? 0 : sample_count + 1; end end // Hanning window processing always @(posedge clk) begin if (rst) begin coeff_count <= 0; sample_out <= 0; sample_out_valid <= 0; end else begin if (sample_count == N-1 && sample_valid) begin coeff_count <= 0; sample_out_valid <= 1; end else if (coeff_count < N) begin coeff <= hanning_coeff[coeff_count]; sample_out <= (sample_buffer[coeff_count] * coeff) >>> Q; coeff_count <= coeff_count + 1; end else if (coeff_count < NF) begin sample_out <= 0; // Zero-padding coeff_count <= coeff_count + 1; end else begin sample_out_valid <= 0; end end end endmodule ", "module log_module #( parameter Q_M = 15, // Number of fractional bits for Mel filter coefficients parameter Q_L = 11 // Number of fractional bits for logarithm output ) ( input wire clk, input wire rst, input wire signed [31:0] data_in, // INT32 Q30 input wire data_valid, output reg signed [15:0] log_out, // INT16 Q11 output reg log_valid ); localparam INT_BITS = 32 - Q_M; localparam FRAC_BITS = Q_M; reg signed [INT_BITS-1:0] int_part; // INT20 Q0 reg signed [FRAC_BITS-1:0] frac_part; // INT15 Q15 reg signed [FRAC_BITS-1:0] frac_part_shifted; // INT15 Q15 reg [$clog2(FRAC_BITS)-1:0] shift_count; reg [1:0] state; always @(posedge clk) begin if (rst) begin int_part <= 0; frac_part <= 0; frac_part_shifted <= 0; shift_count <= 0; log_out <= 0; log_valid <= 0; state <= 0; end else begin case (state) 0: begin if (data_valid) begin int_part <= data_in[31:FRAC_BITS]; frac_part <= data_in[FRAC_BITS-1:0]; state <= 1; end end 1: begin if (int_part > 0) begin int_part <= int_part >> 1; frac_part_shifted <= frac_part; shift_count <= 0; state <= 2; end else begin log_out <= {{1{frac_part[FRAC_BITS-1]}}, frac_part} >> (FRAC_BITS - Q_L); log_valid <= 1; state <= 0; end end 2: begin if (shift_count < FRAC_BITS) begin if (frac_part_shifted >= (1 << (FRAC_BITS - 1))) begin frac_part_shifted <= (frac_part_shifted << 1) - (1 << FRAC_BITS); log_out <= (log_out << 1) + 1; end else begin frac_part_shifted <= frac_part_shifted << 1; log_out <= log_out << 1; end shift_count <= shift_count + 1; end else begin log_out <= log_out[15:0] + (int_part << (Q_L - $clog2(INT_BITS))); log_valid <= 1; state <= 0; end end endcase end end endmodule ","module pow_module #( parameter Q = 15 // Number of fractional bits in the input data ) ( input wire clk, input wire rst, input wire signed [31:0] data_in, // INT32 Q30 input wire data_valid, output reg signed [31:0] data_out, // INT32 Q30 output reg data_out_valid ); reg signed [31:0] data_reg; // INT32 Q30 always @(posedge clk) begin if (rst) begin data_reg <= 0; data_out <= 0; data_out_valid <= 0; end else begin if (data_valid) begin data_reg <= $signed(data_in) * $signed(data_in) >>> Q; data_out <= data_reg; data_out_valid <= 1; end else begin data_out_valid <= 0; end end end endmodule ", "module mel_filterbank #( parameter Q = 15, // Number of fractional bits in the input data parameter NUM_FILTERS = 40, // Number of Mel filters parameter FILTER_SIZE = 23, // Size of each Mel filter parameter Q_M = 15 // Number of fractional bits for Mel filter coefficients ) ( input wire clk, input wire rst, input wire signed [31:0] data_in, // INT32 Q30 input wire data_valid, output reg signed [31:0] mel_out, // INT32 Q30 output reg mel_valid ); localparam NUM_COEFFS = NUM_FILTERS * FILTER_SIZE; localparam COEFF_WIDTH = 16; reg signed [31:0] periodogram [0:FILTER_SIZE-1]; // INT32 Q30 reg signed [15:0] coeff; // INT16 Q15 reg signed [47:0] accumulator; // INT48 Q45 reg [$clog2(NUM_FILTERS)-1:0] filter_counter; reg [$clog2(FILTER_SIZE)-1:0] coeff_counter; reg [1:0] state; // Mel filter coefficients function signed [COEFF_WIDTH-1:0] mel_coeff; input [$clog2(NUM_FILTERS)-1:0] filter_idx; input [$clog2(FILTER_SIZE)-1:0] coeff_idx; reg

signed [COEFF_WIDTH-1:0] result; begin // Implement the Mel filter coefficient calculation here // based on the filter index and coefficient index if (filter_idx < 10) begin if (coeff_idx < 10) result = 16'h7FFF; else if (coeff_idx < 20) result = 16'h4000; else result = 16'h1000; end else if (filter_idx < 30) begin if (coeff_idx < 5) result = 16'h7FFF; else if (coeff_idx < 15) result = 16'h6000; else result = 16'h2000; end else begin if (coeff_idx < 3) result = 16'h7FFF; else if (coeff_idx < 8) result = 16'h5000; else result = 16'h3000; end mel_coeff = result; end endfunction always @(posedge clk) begin if (rst) begin filter_counter <= 0; coeff_counter <= 0; accumulator <= 0; mel_out <= 0; mel_valid <= 0; state <= 0; end else begin case (state) 0: begin if (data_valid) begin periodogram[filter_counter[4:0]] <= data_in; filter_counter <= filter_counter + 1; if (filter_counter == FILTER_SIZE - 1) begin filter_counter <= 0; state <= 1; end end end 1: begin coeff <= mel_coeff(filter_counter, coeff_counter); accumulator <= accumulator + {{16{periodogram[coeff_counter][15]}}, periodogram[coeff_counter]} * {{16{coeff[15]}}, coeff}; coeff_counter <= coeff_counter + 1; if (coeff_counter == FILTER_SIZE - 1) begin mel_out <= accumulator[31:0] >>> (Q + Q_M); mel_valid <= 1; accumulator <= 0; coeff_counter <= 0; filter_counter <= filter_counter + 1; if (filter_counter == NUM_FILTERS) begin filter_counter <= 0; state <= 0; end end end endcase end end endmodule ","module dct_module #( parameter Q_L = 11, // Number of fractional bits for logarithm output parameter Q_D = 4, // Number of fractional bits for DCT output parameter N = 32 // Size of the DCT input vector ) ( input wire clk, input wire rst, input wire signed [15:0] data_in, // INT16 Q11 input wire data_valid, output reg signed [15:0] dct_out, // INT16 Q4 output reg dct_valid ); localparam COEFF_WIDTH = 16; reg signed [11:0] input_buffer [0:N-1]; // INT12 Q11 reg [$clog2(N)-1:0] input_counter; reg [$clog2(N)-1:0] output_counter; reg signed [COEFF_WIDTH-1:0] coeff; reg signed [Q_L+COEFF_WIDTH-1:0] mult; reg signed [Q_L+COEFF_WIDTH-1:0] accumulator; reg [1:0] state; // Precompute DCT coefficients using integer arithmetic wire signed [COEFF_WIDTH-1:0] sqrt_2_over_N = 16'h3B21; // sqrt(2/32) in Q15 format wire signed [COEFF_WIDTH-1:0] sqrt_1_over_N = 16'h2D41; // sqrt(1/32) in Q15 format reg signed [COEFF_WIDTH-1:0] cos_term; reg [$clog2(N)-1:0] x; always @(*) begin if (output_counter == 0) begin coeff = sqrt_1_over_N; x = 0; cos_term = 16'h7FFF; // Initialize to 1 in Q15 format end else begin x = (2 * input_counter + 1) * output_counter; cos_term = 16'h7FFF; // Initialize to 1 in Q15 format repeat (5) begin cos_term = cos_term - ((cos_term * x * x) >> 14); x = x << 1; end coeff = (sqrt_2_over_N * cos_term) >> 15; end end always @(posedge clk) begin if (rst) begin input_counter <= 0; output_counter <= 0; accumulator <= 0; dct_out <= 0; dct_valid <= 0; state <= 0; end else begin case (state) 0: begin if (data_valid) begin input_buffer[input_counter] <= data_in[11:0]; input_counter <= input_counter + 1; if ({{27{1'b0}}, input_counter} == N - 1) begin input_counter <= 0; state <= 1; end end end 1: begin mult <= input_buffer[input_counter] * coeff; accumulator <= accumulator + mult; input_counter <= input_counter + 1; if (input_counter == N - 1) begin dct_out <= accumulator[26:11] >>> (Q_L + COEFF_WIDTH - Q_D); dct_valid <= 1; accumulator <= 0; output_counter <= output_counter + 1; input_counter <= 0; if ({{27{1'b0}}, output_counter} == N - 1) begin output_counter <= 0; state <= 0; end end end endcase end end endmodule ". I receive the following error "Warning: Resizing cell port mfcc_accel.dct.dct_out from 640 bits to 16 bits.". Help me solve it .

ChatGPT

To resolve the error regarding the mismatch in bit widths between **mfcc_accel.dct.dct_out** and the **mfcc_feature** signal, we need to ensure that the output of the DCT module correctly generates 40 features each of 16 bits, and **mfcc_accel** correctly handles these features.

Here is a step-by-step breakdown of the required changes:

1. **Refactor the dct_module to output 40 features of 16 bits each**: The DCT module needs to output 40 features instead of a single 16-bit value. We need to change its implementation to process multiple features and output them in parallel or sequentially as appropriate.

2. **Update the mfcc_accel to handle the 640-bit mfcc_feature output**: Ensure that the **mfcc_accel** module correctly aggregates the 16-bit outputs of the DCT module into a 640-bit output.

**Refactor dct_module**

The DCT module will be refactored to output 40 features of 16 bits each. This can be achieved by accumulating the DCT results and storing them in an array.

verilog

Copy code

module dct_module #( parameter Q_L = 11, // Number of fractional bits for logarithm output parameter Q_D = 4, // Number of fractional bits for DCT output parameter N = 32 // Size of the DCT input vector ) ( input wire clk, input wire rst, input wire signed [15:0] data_in, // INT16 Q11 input wire data_valid, output reg [639:0] dct_out, // 40 features * 16 bits = 640 bits output reg dct_valid ); localparam COEFF_WIDTH = 16; reg signed [15:0] input_buffer [0:N-1]; // INT16 Q11 reg [$clog2(N)-1:0] input_counter; reg [5:0] output_counter; reg signed [COEFF_WIDTH-1:0] coeff; reg signed [Q_L+COEFF_WIDTH-1:0] mult; reg signed [Q_L+COEFF_WIDTH-1:0] accumulator; reg [1:0] state; // Precompute DCT coefficients using integer arithmetic wire signed [COEFF_WIDTH-1:0] sqrt_2_over_N = 16'h3B21; // sqrt(2/32) in Q15 format wire signed [COEFF_WIDTH-1:0] sqrt_1_over_N = 16'h2D41; // sqrt(1/32) in Q15 format reg signed [COEFF_WIDTH-1:0] cos_term; reg [$clog2(N)-1:0] x; always @(*) begin if (output_counter == 0) begin coeff = sqrt_1_over_N; x = 0; cos_term = 16'h7FFF; // Initialize to 1 in Q15 format end else begin x = (2 * input_counter + 1) * output_counter; cos_term = 16'h7FFF; // Initialize to 1 in Q15 format repeat (5) begin cos_term = cos_term - ((cos_term * x * x) >> 14); x = x << 1; end coeff = (sqrt_2_over_N * cos_term) >> 15; end end always @(posedge clk) begin if (rst) begin input_counter <= 0; output_counter <= 0; accumulator <= 0; dct_out <= 0; dct_valid <= 0; state <= 0; end else begin case (state) 0: begin if (data_valid) begin input_buffer[input_counter] <= data_in; input_counter <= input_counter + 1; if (input_counter == N - 1) begin input_counter <= 0; state <= 1; end end end 1: begin mult <= input_buffer[input_counter] * coeff; accumulator <= accumulator + mult; input_counter <= input_counter + 1; if (input_counter == N - 1) begin dct_out[output_counter*16 +: 16] <= accumulator[26:11] >>> (Q_L + COEFF_WIDTH - Q_D); accumulator <= 0; output_counter <= output_counter + 1; input_counter <= 0; if (output_counter == 39) begin output_counter <= 0; dct_valid <= 1; state <= 0; end end end endcase end end endmodule

**Update mfcc_accel**

Ensure that the **mfcc_accel** module correctly handles the 640-bit output from the **dct_module**.

verilog

Copy code

module mfcc_accel ( input wire clk, input wire rst, input wire [15:0] audio_sample, input wire sample_valid, output reg [639:0] mfcc_feature, // 40 features * 16 bits = 640 bits output reg mfcc_valid ); // Signals for loading weights reg load_weights; wire [15:0] hanning_real_out; wire hanning_real_valid; // Instantiate the Hanning window real module hanning_window_real

hanning_real ( .clk(clk), .rst(rst), .sample_in(audio_sample), .sample_valid(sample_valid), .sample_out(hanning_real_out), .sample_out_valid(hanning_real_valid) ); // Instantiate the Hanning window imaginary module wire [15:0] hanning_imag_out; wire hanning_imag_valid; hanning_window_imag hanning_imag ( .clk(clk), .rst(rst), .sample_in(audio_sample), .sample_valid(sample_valid), .sample_out(hanning_imag_out), .sample_out_valid(hanning_imag_valid) ); // Instantiate the periodogram module wire [31:0] periodogram_out; wire periodogram_valid; periodogram_squared periodogram ( .clk(clk), .rst(rst), .sample_in_real(hanning_real_out), .sample_in_imag(hanning_imag_out), .sample_valid(hanning_real_valid), // Use the valid signal from the real Hanning window .periodogram_out(periodogram_out), .periodogram_valid(periodogram_valid) ); // Instantiate the pow_module wire [31:0] pow_out; wire pow_valid; pow_module pow ( .clk(clk), .rst(rst), .data_in(periodogram_out), .data_valid(periodogram_valid), .data_out(pow_out), .data_out_valid(pow_valid) ); // Instantiate the mel_filterbank module wire [31:0] mel_out; wire mel_valid; mel_filterbank mel ( .clk(clk), .rst(rst), .data_in(pow_out), .data_valid(pow_valid), .mel_out(mel_out), .mel_valid(mel_valid) ); // Instantiate the log_module wire [15:0] log_out; wire log_valid; log_module log ( .clk(clk), .rst(rst), .data_in(mel_out), .data_valid(mel_valid), .log_out(log_out), .log_valid(log_valid) ); // Instantiate the dct_module wire [639:0] dct_out; wire dct_valid; dct_module dct ( .clk(clk), .rst(rst), .data_in(log_out), .data_valid(log_valid), .dct_out(dct_out), .dct_valid(dct_valid ); always @(posedge clk or posedge rst) begin if (rst) begin mfcc_feature <= 0; mfcc_valid <= 0; end else begin mfcc_feature <= dct_out; mfcc_valid <= dct_valid; end end endmodule

**Update cnn_kws_accel**

Now that the **mfcc_accel** correctly produces a 640-bit output for the MFCC features, ensure that **cnn_kws_accel** correctly handles the 640-bit **mfcc_feature** output from **mfcc_accel**.

verilog

Copy code

module cnn_kws_accel #( parameter NUM_KEYWORDS = 10, parameter MFCC_FEATURES = 40, parameter ACTIV_BITS = 16, // Each MFCC feature is 16 bits parameter FC1_INPUT_SIZE = (MFCC_FEATURES/2) * (CONV2_NUM_FILTERS), parameter FC1_OUTPUT_SIZE = 64, parameter FC2_INPUT_SIZE = 64, parameter FC2_OUTPUT_SIZE = NUM_KEYWORDS, parameter CONV1_KERNEL_SIZE = 3, parameter CONV1_NUM_FILTERS = 8, parameter CONV2_KERNEL_SIZE = 3, parameter CONV2_NUM_FILTERS = 8 )( input wire clk, input wire rst_n, input wire [15:0] audio_in, input wire audio_valid, output reg [NUM_KEYWORDS-1:0] kws_result, output reg kws_valid, input wire [7:0] frame_size, input wire [7:0] frame_overlap, input wire [7:0] num_mfcc_coeffs, input wire [7:0] goertzel_coefs, output wire goertzel_coefs_start, output wire goertzel_coefs_valid, output wire goertzel_coefs_done, input wire [CONV1_NUM_FILTERS*CONV1_KERNEL_SIZE*CONV1_KERNEL_SIZE*ACTIV_BITS-1:0] conv1_weights, input wire [CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_biases, input wire conv1_load_weights, input wire conv1_load_biases, input wire [CONV2_NUM_FILTERS*CONV1_NUM_FILTERS*CONV2_KERNEL_SIZE*CONV2_KERNEL_SIZE*ACTIV_BITS-1:0] conv2_weights, input wire [CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_biases, input wire conv2_load_weights, input wire conv2_load_biases, input wire [FC1_OUTPUT_SIZE*FC1_INPUT_SIZE*ACTIV_BITS-1:0] fc1_weights, input wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_biases, input wire fc1_load_weights, input wire fc1_load_biases, input wire [FC2_OUTPUT_SIZE*FC2_INPUT_SIZE*ACTIV_BITS-1:0] fc2_weights,

input wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_biases, input wire fc2_load_weights, input wire fc2_load_biases ); // MFCC module signals wire [639:0] mfcc_out; // 40 features * 16 bits wire mfcc_valid; // CNN-KWS layers wire [MFCC_FEATURES*CONV1_NUM_FILTERS*ACTIV_BITS-1:0] conv1_out; wire conv1_valid; wire [MFCC_FEATURES*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] conv2_out; wire conv2_valid; wire [(MFCC_FEATURES/2)*CONV2_NUM_FILTERS*ACTIV_BITS-1:0] maxpool_out; wire maxpool_valid; wire [FC1_OUTPUT_SIZE*ACTIV_BITS-1:0] fc1_out; wire fc1_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] fc2_out; wire fc2_valid; wire [FC2_OUTPUT_SIZE*ACTIV_BITS-1:0] softmax_out; wire softmax_valid; // MFCC module instantiation mfcc_accel mfcc ( .clk(clk), .rst(rst_n), // Adjusted to rst_n .audio_sample(audio_in), // Adjusted to audio_sample .sample_valid(audio_valid), // Adjusted to sample_valid .mfcc_feature(mfcc_out), // Adjusted to mfcc_feature .mfcc_valid(mfcc_valid) // Adjusted to mfcc_valid ); // Convolutional layer 1 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(1), .KERNEL_SIZE(CONV1_KERNEL_SIZE), .NUM_FILTERS(CONV1_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv1 ( .clk(clk), .rst_n(rst_n), .data_in(mfcc_out), .data_valid(mfcc_valid), .data_out(conv1_out), .data_out_valid(conv1_valid), .weights_in(conv1_weights), .biases_in(conv1_biases), .load_weights(conv1_load_weights), .load_biases(conv1_load_biases) ); // Convolutional layer 2 conv2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV1_NUM_FILTERS), .KERNEL_SIZE(CONV2_KERNEL_SIZE), .NUM_FILTERS(CONV2_NUM_FILTERS), .PADDING(1), .ACTIV_BITS(ACTIV_BITS) ) conv2 ( .clk(clk), .rst_n(rst_n), .data_in(conv1_out), .data_valid(conv1_valid), .data_out(conv2_out), .data_out_valid(conv2_valid), .weights_in(conv2_weights), .biases_in(conv2_biases), .load_weights(conv2_load_weights), .load_biases(conv2_load_biases) ); maxpool2d #( .INPUT_WIDTH(MFCC_FEATURES), .INPUT_HEIGHT(1), .INPUT_CHANNELS(CONV2_NUM_FILTERS), .KERNEL_SIZE(2), .STRIDE(2), .ACTIV_BITS(ACTIV_BITS) ) maxpool ( .clk(clk), .rst_n(rst_n), .data_in(conv2_out), .data_valid(conv2_valid), .data_out(maxpool_out), .data_out_valid(maxpool_valid) ); fully_connected #( .INPUT_SIZE(FC1_INPUT_SIZE), .OUTPUT_SIZE(FC1_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc1 ( .clk(clk), .rst_n(rst_n), .data_in(maxpool_out), .data_valid(maxpool_valid), .data_out(fc1_out), .data_out_valid(fc1_valid), .weights_in(fc1_weights), .biases_in(fc1_biases), .load_weights(fc1_load_weights), .load_biases(fc1_load_biases) ); // Fully connected layer 2 (output layer) fully_connected #( .INPUT_SIZE(FC2_INPUT_SIZE), .OUTPUT_SIZE(FC2_OUTPUT_SIZE), .ACTIV_BITS(ACTIV_BITS) ) fc2 ( .clk(clk), .rst_n(rst_n), .data_in(fc1_out), .data_valid(fc1_valid), .data_out(fc2_out), .data_out_valid(fc2_valid), .weights_in(fc2_weights), .biases_in(fc2_biases), .load_weights(fc2_load_weights), .load_biases(fc2_load_biases) ); // Softmax layer softmax #( .INPUT_SIZE(NUM_KEYWORDS), .ACTIV_BITS(ACTIV_BITS) ) softmax ( .clk(clk), .rst_n(rst_n), .data_in(fc2_out), .data_valid(fc2_valid), .data_out(softmax_out), .data_out_valid(softmax_valid) ); // Output assignment always @(posedge clk or negedge rst_n) begin if (!rst_n) begin kws_result <= 'b0; kws_valid <= 1'b0; end else begin kws_result <= softmax_out[NUM_KEYWORDS-1:0]; kws_valid <= softmax_valid; end end endmodule

**Summary of Changes:**

1. **dct_module**:

   - Refactored to produce 40 features of 16 bits each, concatenated into a 640-bit output.

2. **mfcc_accel**:

   - Updated to handle 640-bit **mfcc_feature** output from **dct_module**.

3. **cnn_kws_accel**:

- Adjusted to handle 640-bit **mfcc_feature** input from **mfcc_accel**.

By making these changes, the warning and bit-width mismatch issues should be resolved, and the system should correctly process the MFCC features through the CNN layers to perform keyword spotting.