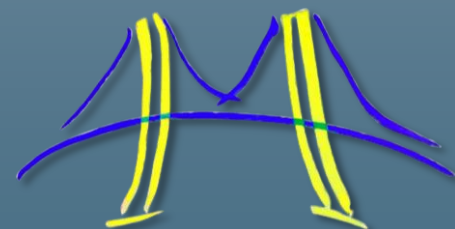
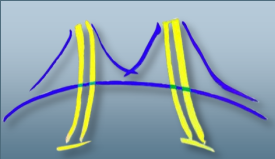


PRODUCTIVE GMM TRAINING WITH SEJITS FOR SPEAKER DIARIZATION

Katya Gonina, Henry Cook, Shoiab Kamil,
Gerald Friedland, Armando Fox, David Patterson

ParLab Retreat, June 2, 2011






The Meeting Diarist

Applet Viewer: jokeomat.Jokeomat

Video



Filter

Filter by Keyword:


OK

Filter by Person:

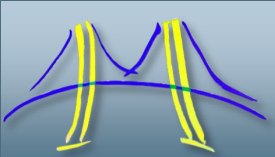
- ☒ Jerry
- ☒ George
- ☐ female
- ☒ Elaine
- ☒ male

Navigation

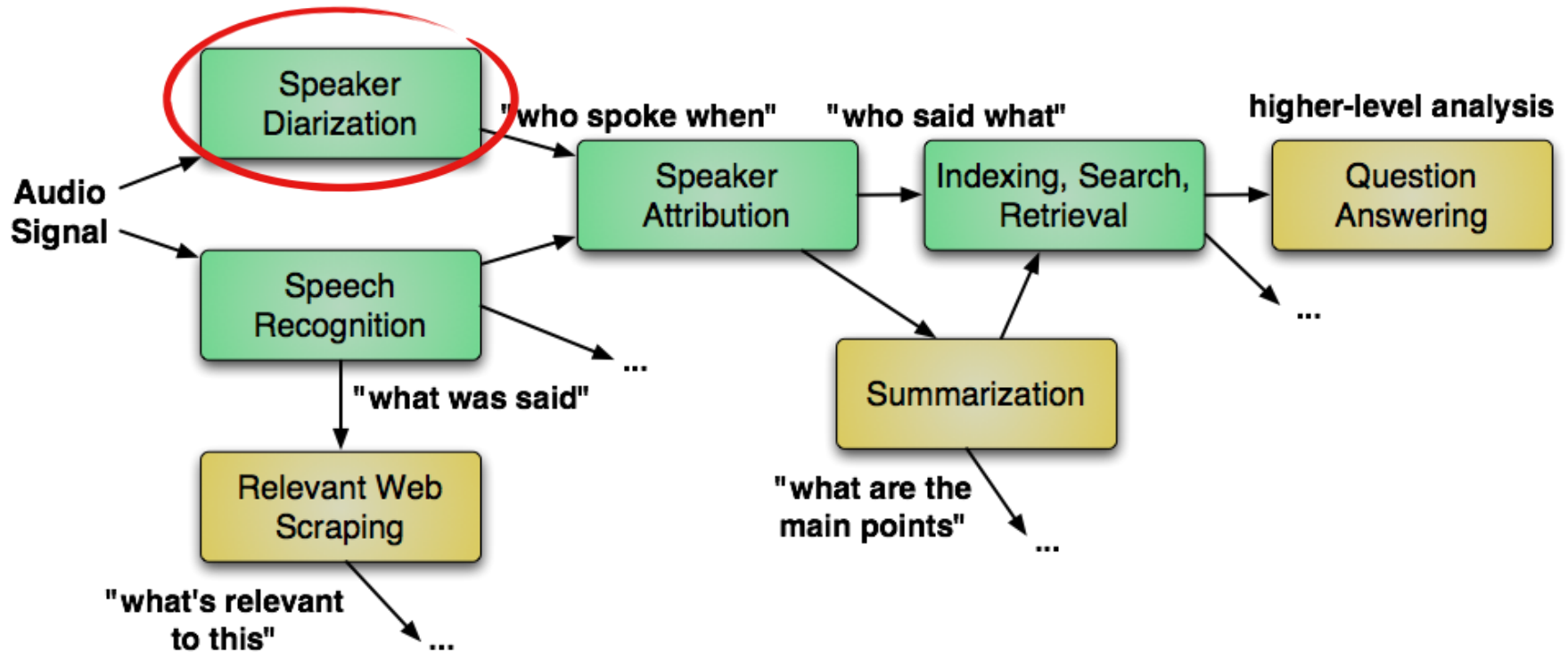
Scenes Top-5 Punchlines Punchlines Dialog

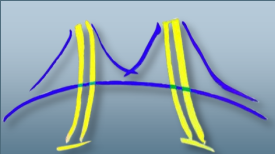


Applet started.



Components of the Meeting Diarist





Speaker Diarization

Audio track:



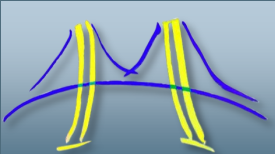
Segmentation:



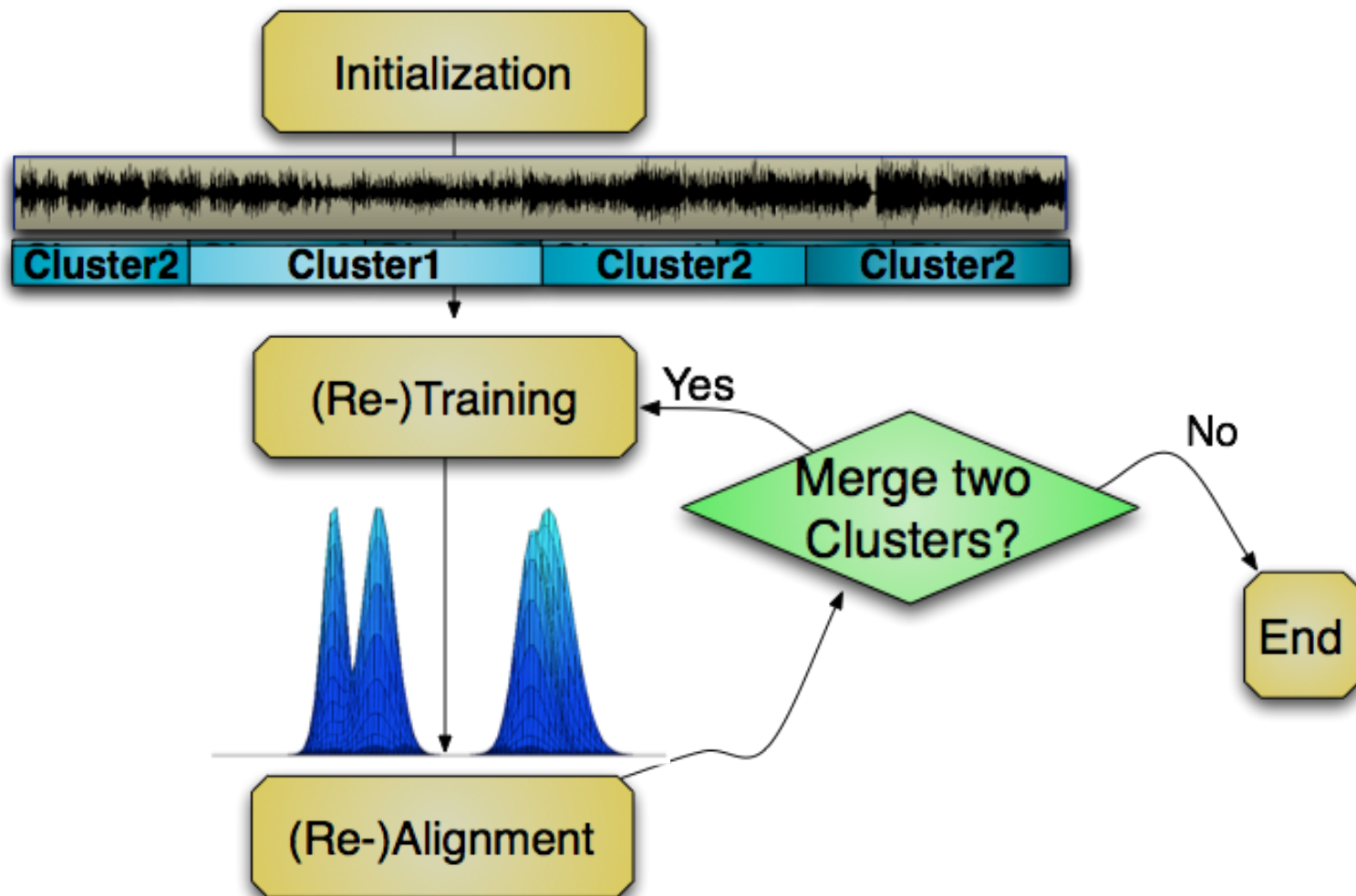
Clustering:

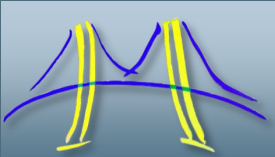


Estimate “who spoke when” with no prior knowledge of speakers, #of speakers, words, or language spoken.



Speaker Diarization: Core Algorithm

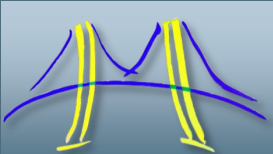




Parallelization of Diarization

- Five versions (so far):
 - ❖ Initial code (2006): 0.333 x realtime (i.e., 1h audio = 3h processing)
 - ❖ Serially optimized (2008): 1.5 x realtime
 - ❖ Parlab retreat summer 2010: Multicore+GPU parallelization: 14.3 x realtime
 - ❖ Parlab retreat winter 2011: GPU-only parallelization 250 x realtime (i.e., 1h audio = 14.4sec processing)
 - ❖ -> Offline = online!
 - ❖ **Parlab retreat summer 2011: SEJITized! [1]**

[1] H. Cook, E. Gonina, S. Kamil, G. Friedland, D. Patterson, A. Fox. CUDA-level Performance with Python-level Productivity for Gaussian Mixture Model Applications. USENIX HotPar Workshop, 2011.



Speaker Diarization in Python

```
def AHC(self):
```

```
# Get the events, divide them into an initial k clusters and train each GMM on a cluster
per_cluster = self.N/self.init_num_clusters
init_training = zip(self.gmm_list,np.vsplit(self.X, range(per_cluster, self.N, per_cluster)))
for g, x in init_training:
    g.train(x)
```

```
# Perform hierarchical agglomeration based on BIC scores
best_BIC_score = 1.0
while (best_BIC_score > 0 and len(self.gmm_list) > 1):
```

```
    num_clusters = len(self.gmm_list)
```

```
    # Resegment data based on likelihood scoring
    likelihoods = self.gmm_list[0].score(self.X)
    for g in self.gmm_list[1:]:
        likelihoods = np.column_stack((likelihoods, g.score (self.X)))
    most_likely = likelihoods.argmax(axis=1)
```

```
    # Across 2.5 secs of observations, vote on which cluster they should be associated with
    split_events = split_events_based_on_votes(most_likely, self.X)
```

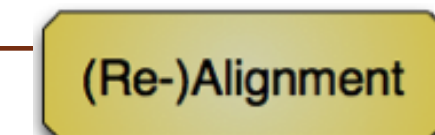
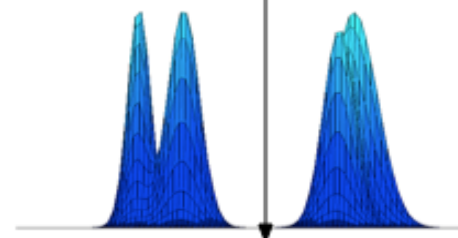
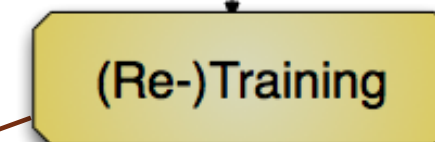
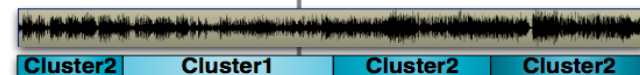
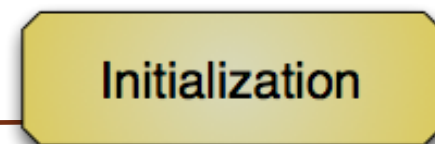
```
    for g, data in split_events:
        g.train(data)
```

```
    # Score all pairs of GMMs using BIC
    best_merged_gmm = None
    best_BIC_score = 0.0
    merged_tuple = None
```

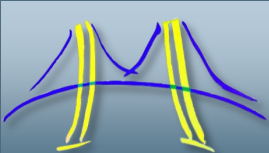
```
    for gmm1idx in range(len(iter_bic_list)):
        for gmm2idx in range(gmm1idx+1, len(iter_bic_list)):
            g1, d1 = iter_bic_list[gmm1idx]
            g2, d2 = iter_bic_list[gmm2idx]
            score = 0.0
            new_gmm, score = compute_distance_BIC(g1, g2, np.concatenate((d1, d2)))
            if score > best_BIC_score:
                best_merged_gmm = new_gmm
                merged_tuple = (g1, g2)
                best_BIC_score = score
```

```
    # Merge the winning candidate pair
```

```
    if best_BIC_score > 0.0:
        self.gmm_list.remove(merged_tuple[0])
        self.gmm_list.remove(merged_tuple[1])
        self.gmm_list.append(best_merged_gmm)
```



Yes



Speaker Diarization in Python

```
def AHC(self):
```

```
# Get per-cluster GMMs and train each GMM on a cluster
init new_gmm_list(M,D)
for g in gmm_list:
    g.train(x)
```

```
# Perform hierarchical agglomeration based on BIC scores
best_BIC_score = 1.0
while (best_BIC_score > 0 and len(self.gmm_list) > 1):
```

```
    num_clusters = len(self.gmm_list)
```

```
# Resegment data based on likelihood scoring
likelihoods = self.gmm_list[0].score(self.X)
for g in self.gmm_list[1:]:
    likelihoods = np.column_stack((likelihoods, g.score(self.X)))
most_likely = likelihoods.argmax(axis=1)
```

```
# Across 2.5 secs of observations, vote on which cluster they should be associated with
split_idx = self.X[most_likely].on_votes(most_likely, self.X)
```

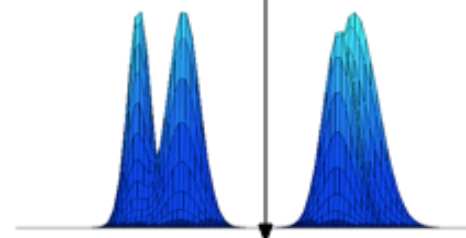
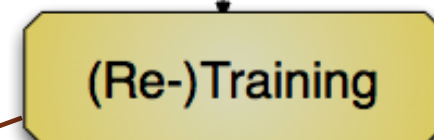
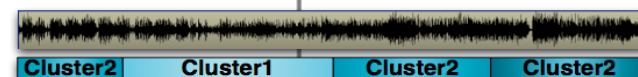
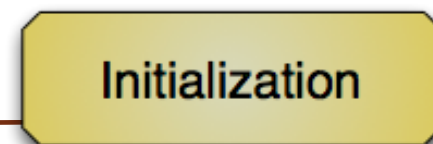
```
for g, x in zip(self.gmm_list, self.X):
    g.train(x)
```

```
# Score all pairs of GMMs using BIC
best_merged_gmm = None
best_BIC_score = 0.0
merged_tuple = None
```

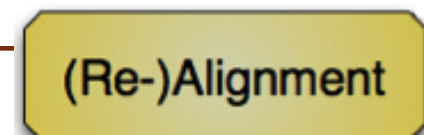
```
for gmm1idx in range(len(iter_bic_list)):
    for gmm2idx in range(gmm1idx+1, len(iter_bic_list)):
        g1, d1 = iter_bic_list[gmm1idx]
        g2, d2 = iter_bic_list[gmm2idx]
        score = 0.0
        new_gmm, score = combine(g1, g2)
        if score > best_BIC_score:
            best_merged_gmm = new_gmm
            merged_tuple = (g1, g2)
            best_BIC_score = score
```

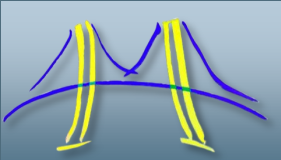
```
# Merge the winning candidate pair
```

```
if best_BIC_score > 0.0:
    self.gmm_list.remove(merged_tuple[0])
    self.gmm_list.remove(merged_tuple[1])
    self.gmm_list.append(best_merged_gmm)
```

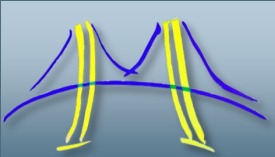


Yes



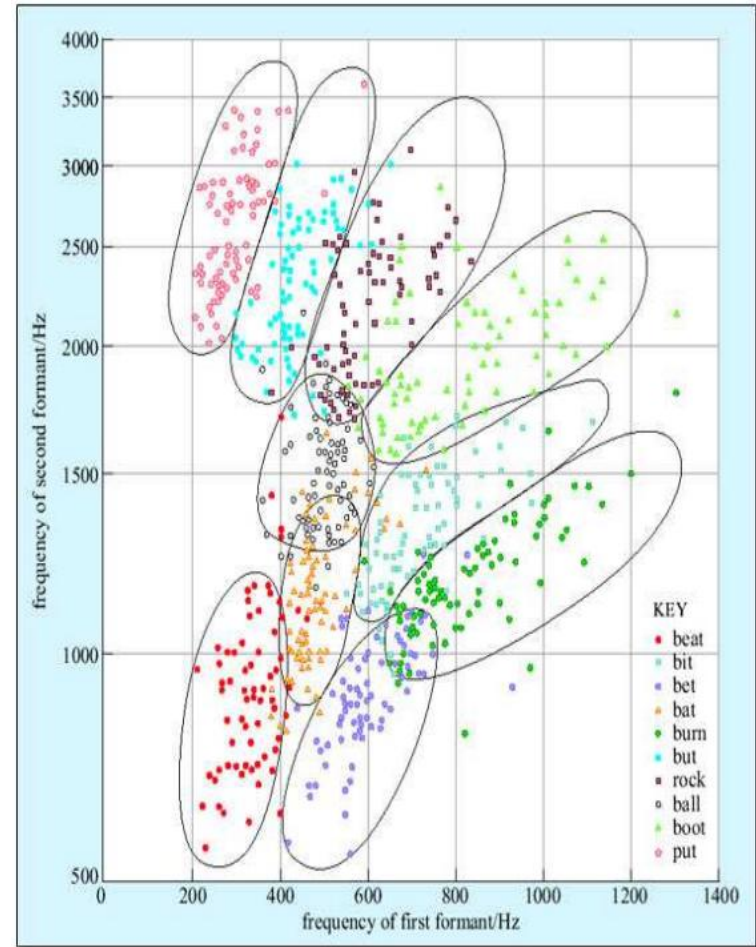


Gaussian Mixture Models & Training

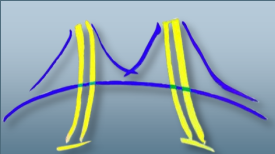


Clustering with Gaussian Mixture Models

- ▶ GMM - probabilistic model for clustering (audio) data
 - ▶ Assumes the distribution of observations follows a set (mixture) of multidimensional Gaussian distributions
 - ▶ Each Gaussian in the mixture has a mean (μ) and a covariance (Σ) parameters
 - ▶ Gaussians in the mixture are weighted with weight D

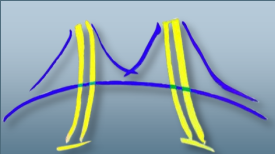


Example GMM in two dimensions
(Source: Dan Klein, UCB)



GMM Training using EM Algorithm

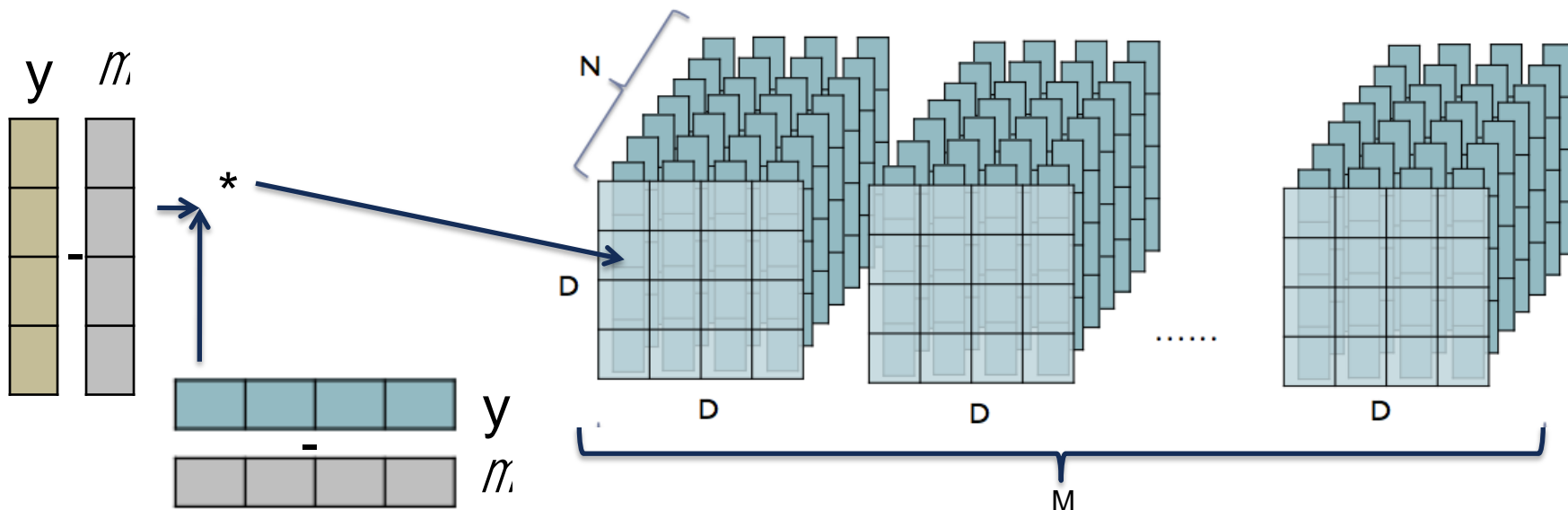
- Given a set of observations/events – find the maximum likelihood estimates of the set of Gaussian Mixture parameters (π, μ, Σ) and classify observations
- Expectation Maximization (EM) Algorithm
 - E step
 - Compute probabilities of events given model parameters
 - M step
 - Compute model parameters given probabilities
 - weights, mean, **covariance matrix**
 - Iterate until convergence
- Covariance matrix – most computationally intensive step

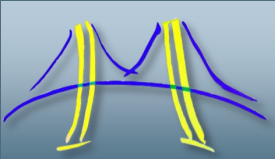


Covariance Matrix Computation

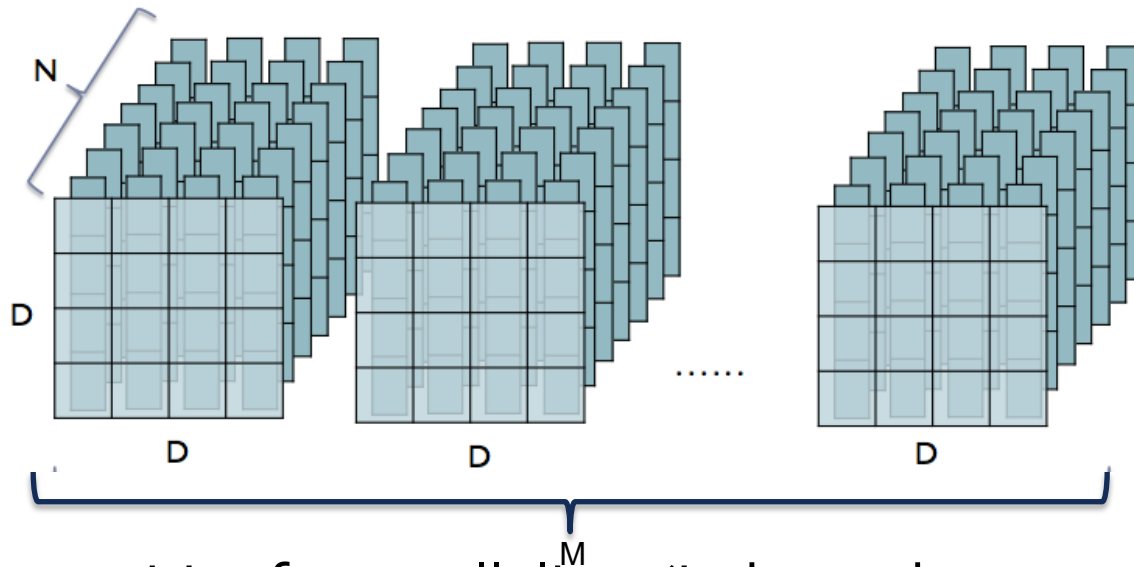
- N – number of feature vectors, ~10K-100K
- D – feature vector dimension (19 for speaker diarization), ~10-100
- M – number of Gaussian components, ~1-128
- Matrix is symmetric – only compute the lower DxD/2 cells

$$\bar{R}_k = \frac{1}{\bar{N}_k} \sum_{n=1}^N (y_n - \bar{\mu}_k)(y_n - \bar{\mu}_k)^t p_{x_n|y_n}(k|y_n, \theta^{(i)})$$

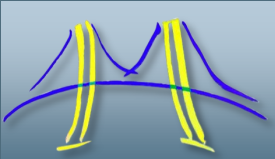




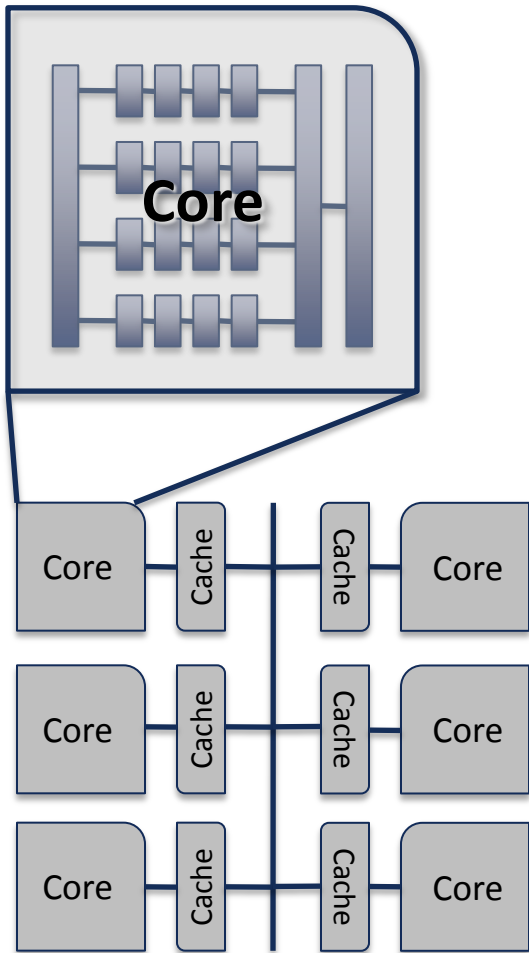
Covariance Matrix Computation



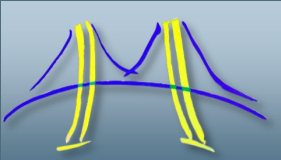
- Opportunities for parallelism (independent computations):
 - Each component's covariance matrix (M)
 - Each cell in a covariance matrix ($D \times D / 2$)
 - Each event's contribution to a cell in a covar matrix (N)
- -> **Multiple code variants** to perform the same computation in different ways



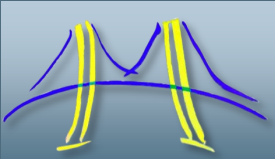
Manycore Parallel Platform



- Two levels of parallelism:
 - Work-groups – parallelized across cores (CUDA threadBlock)
 - Work-groups' work-items – executed on a single core, utilizing within-core parallelism (CUDA thread)
- Per-core local memory

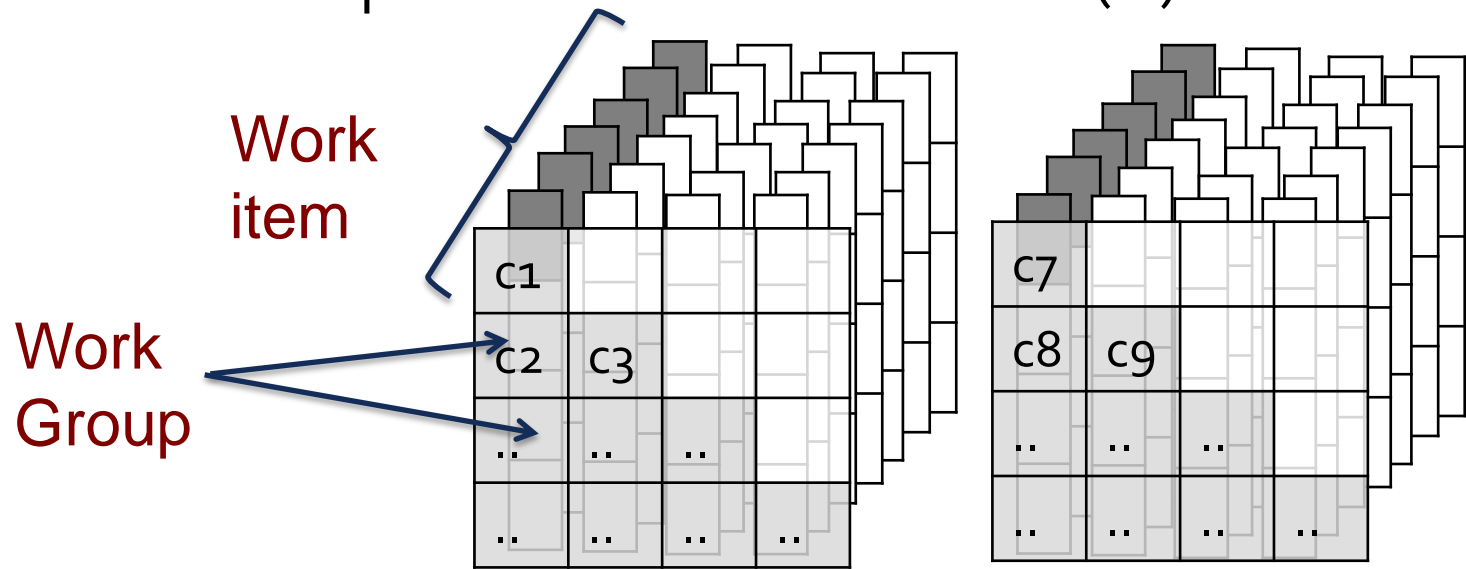


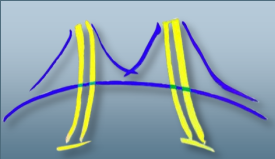
Code Variants



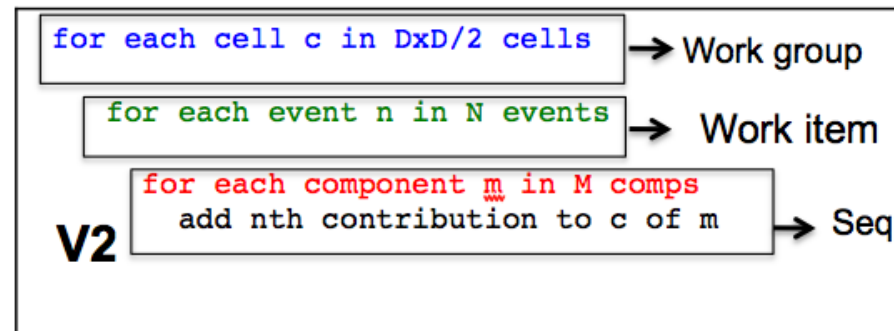
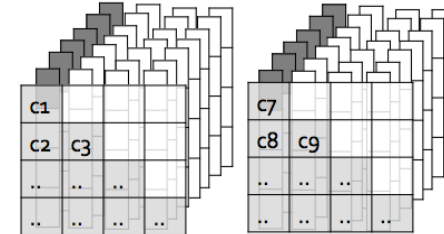
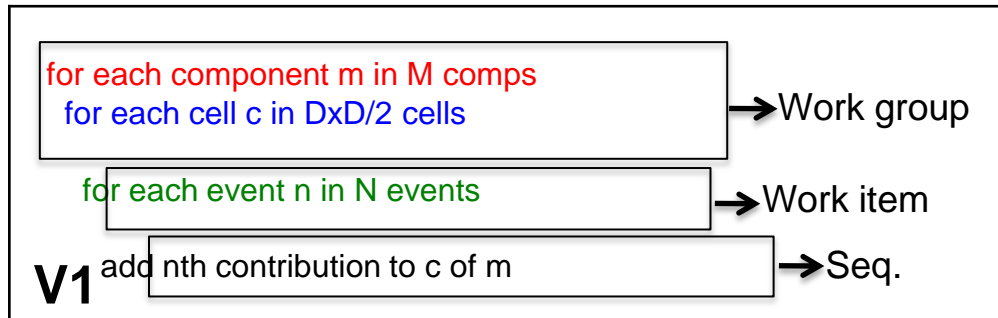
Code Variants - Example

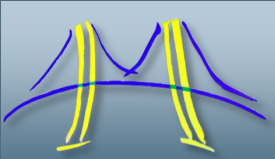
- Code variant 1:
 - 2D grid of work groups $M \times D \times D/2$
 - Each work group is responsible for computing **one cell** in the covariance matrix for **one component**
 - Work item parallelization over events (N)





Covariance Matrix Computation – Code Variants





Covariance Matrix Computation – Code Variants Summary

V1

```

for each component m in M comps
  for each cell c in Dx D/2 cells
    for each event n in N events
      add nth contribution to c of m
  
```

→ Work group
→ Work item
→ Seq

V2

```

for each cell c in Dx D/2 cells
  for each event n in N events
    for each component m in M comps
      add nth contribution to c of m
  
```

→ Work group
→ Work item
→ Seq

V3

```

for each component m in M comps
  for each cell c in Dx D/2 cells
    for each event n in N events
      add nth contribution to c of m
  
```

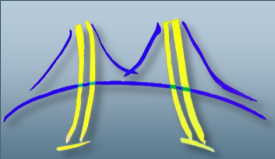
→ Work group
→ Work item
→ Seq

V4

```

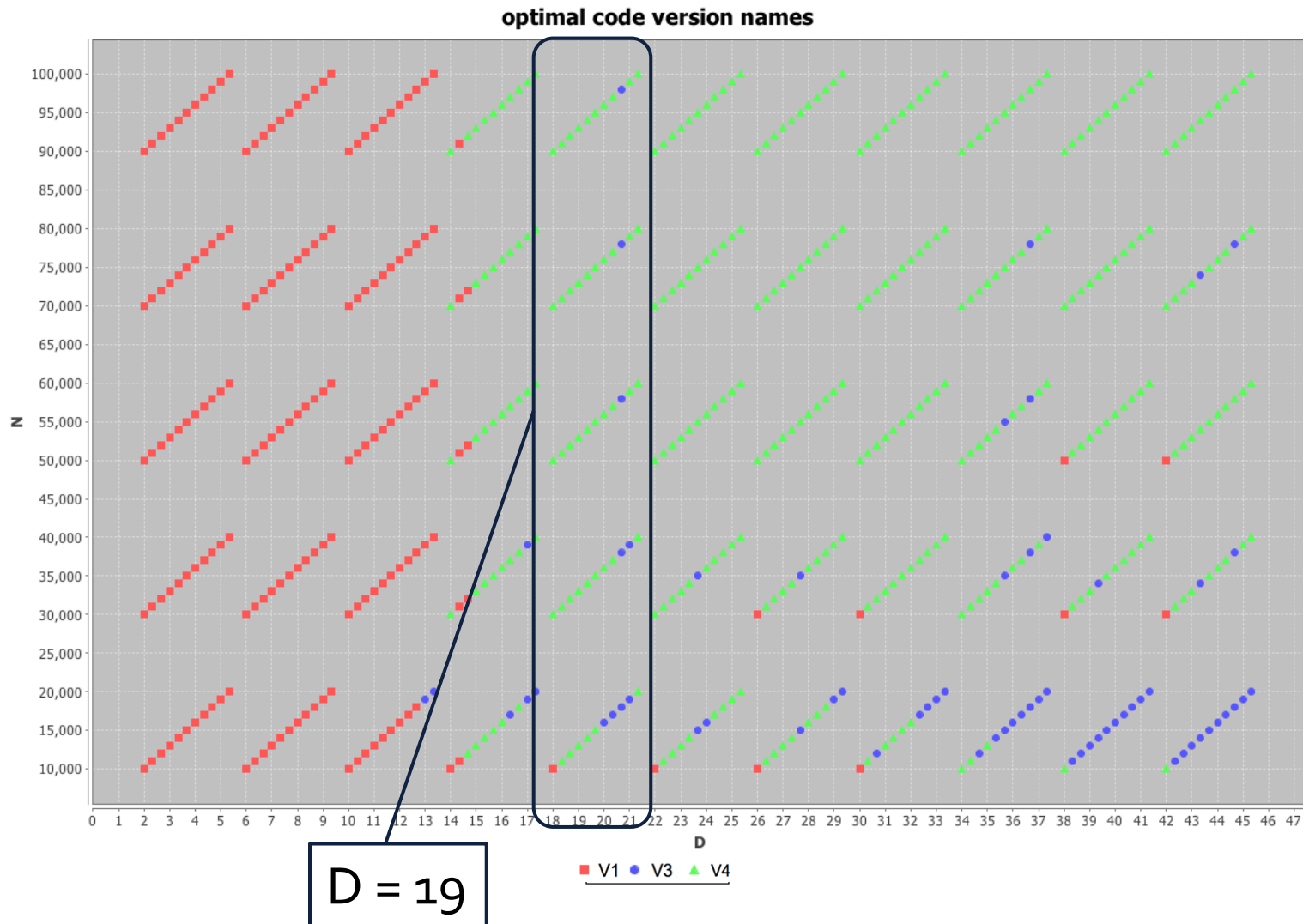
for each block b in B event blocks
  for each component m in M comps
    for each cell c in Dx D/2 cells
      for each event n in N/B events
        add nth contribution to c of m
    for each component m in M comps
      for each block b in B event blocks
        sum partial contributions to m from b
  
```

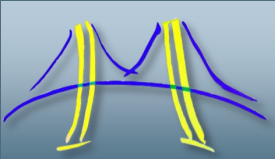
→ Work group
→ Work item
→ Seq
→ Seq



Results – Code Variant Performance

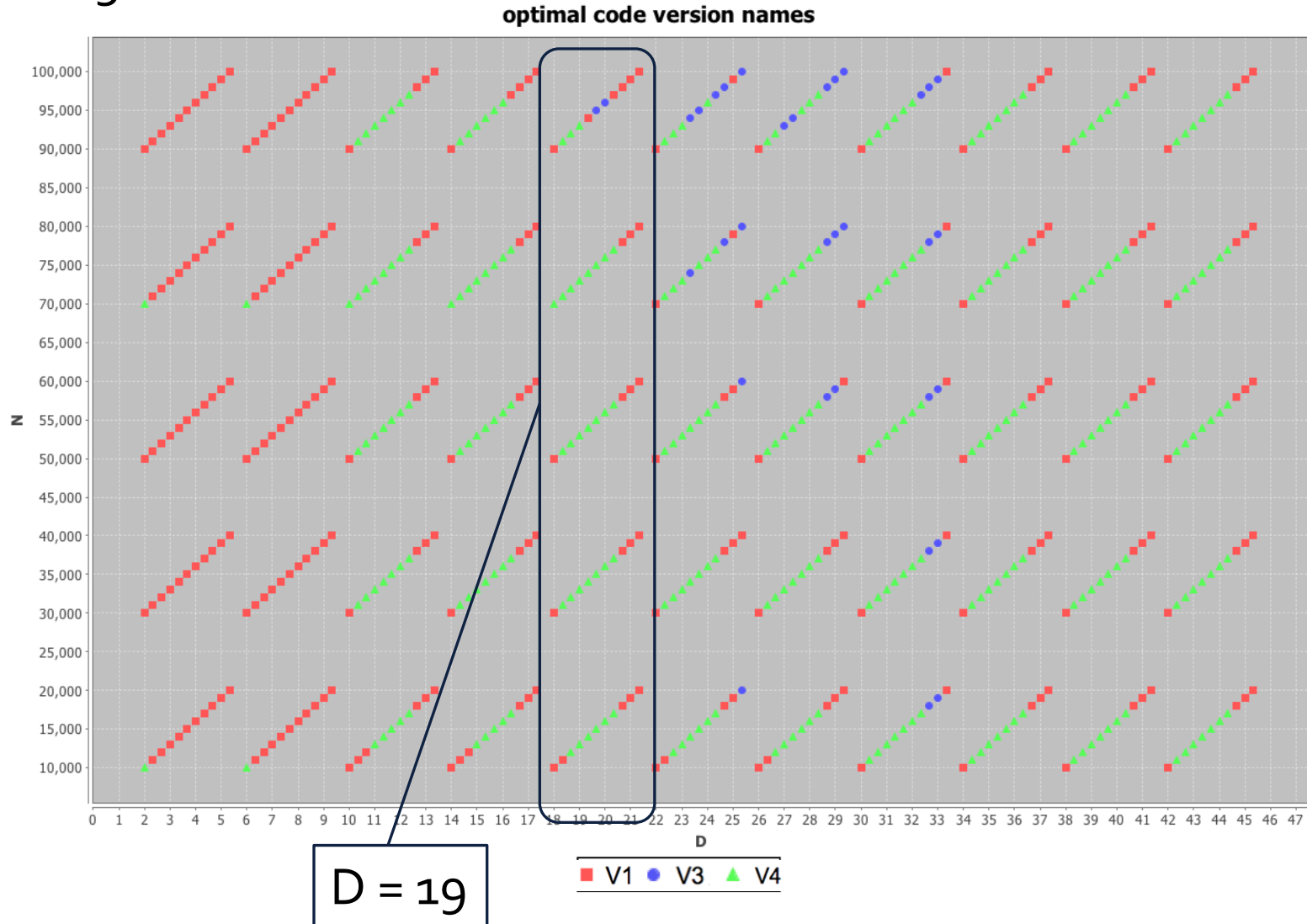
GTX₄₈₀

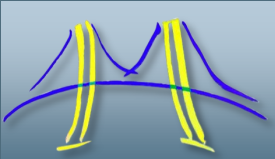




Results – Code Variant Performance

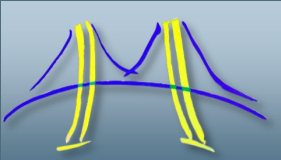
GTX285



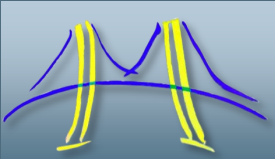


Results - Code Variant Selection

- Using best-performing code variant gave 32% average improvement in matrix computation time compared to always using original hand-coded variant (D: 1 to 36, M: 1 to 128, N: 10K to 150K)
- Performance gap increases with larger problem sizes (75.6% for D=36, M=128, N=500,000)

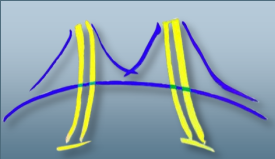


Specialization

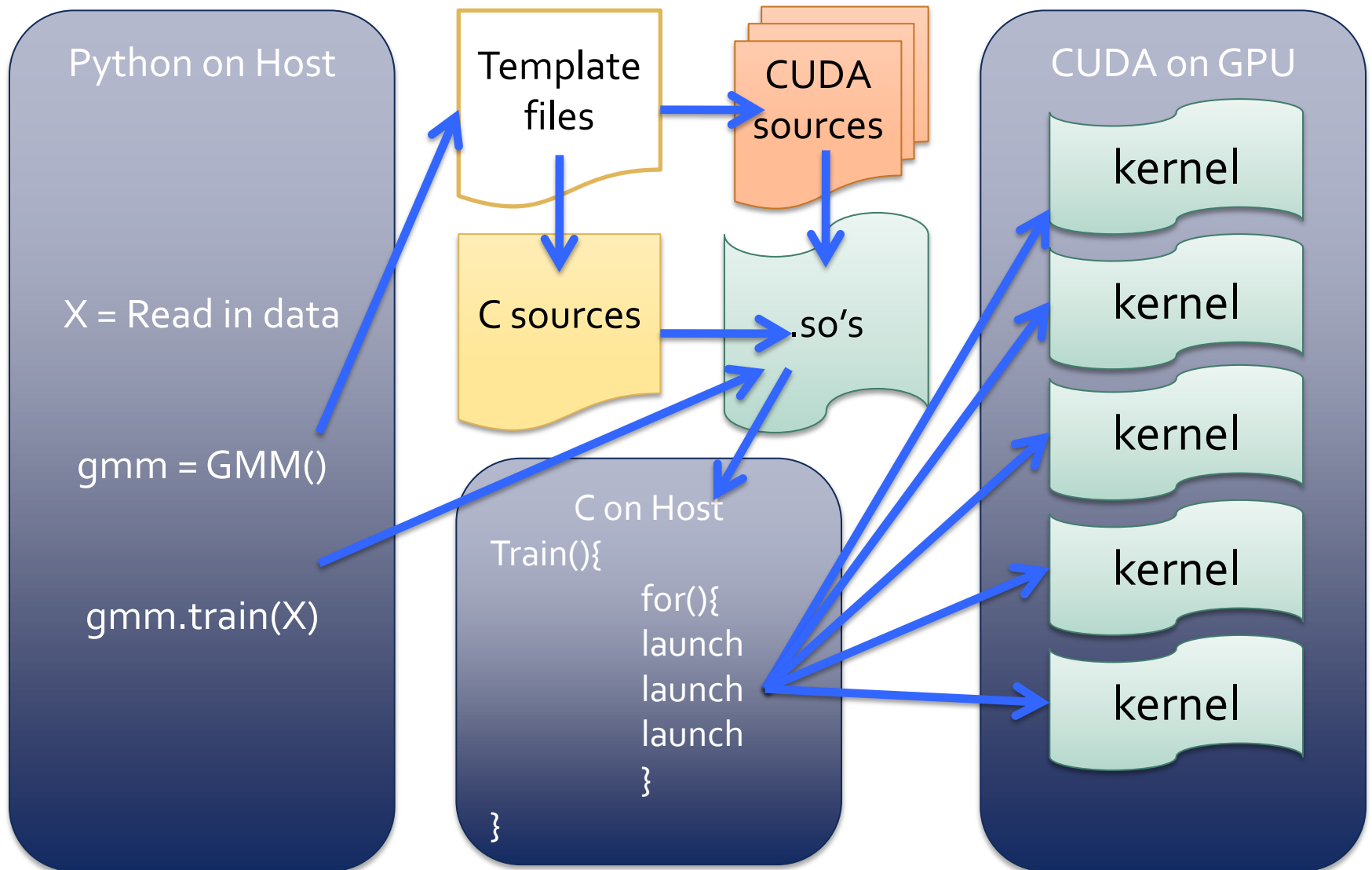


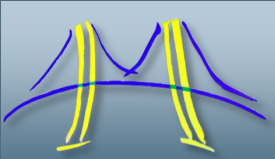
Specialization with ASP

- Given:
 - Problem Dimensions (N, D, M)
 - Platform Parameters (targeting Nvidia GPUs)
 - Core count, shared memory size, SIMD width...
- Automatically select:
 - Optimal code variant
 - Optimal parameters (block size, number of blocks) for that parallelization strategy



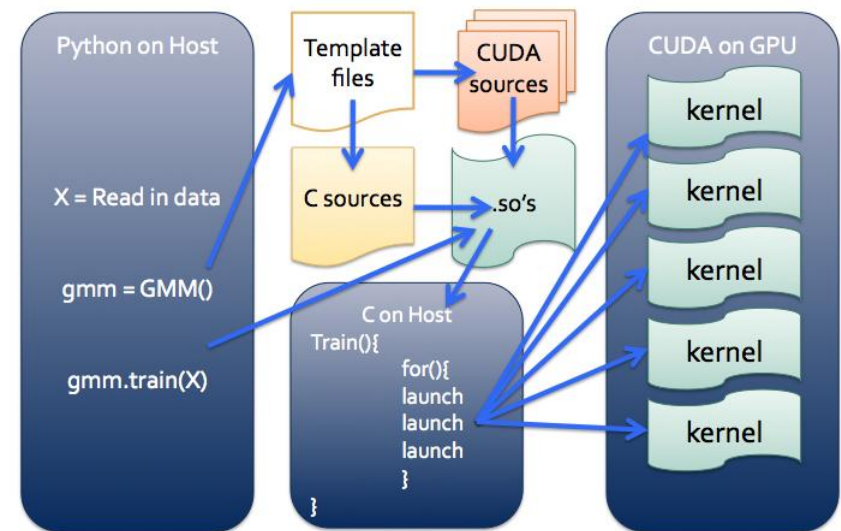
SEJITS Framework: Overview

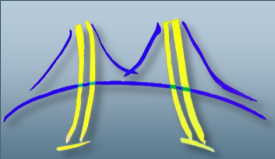




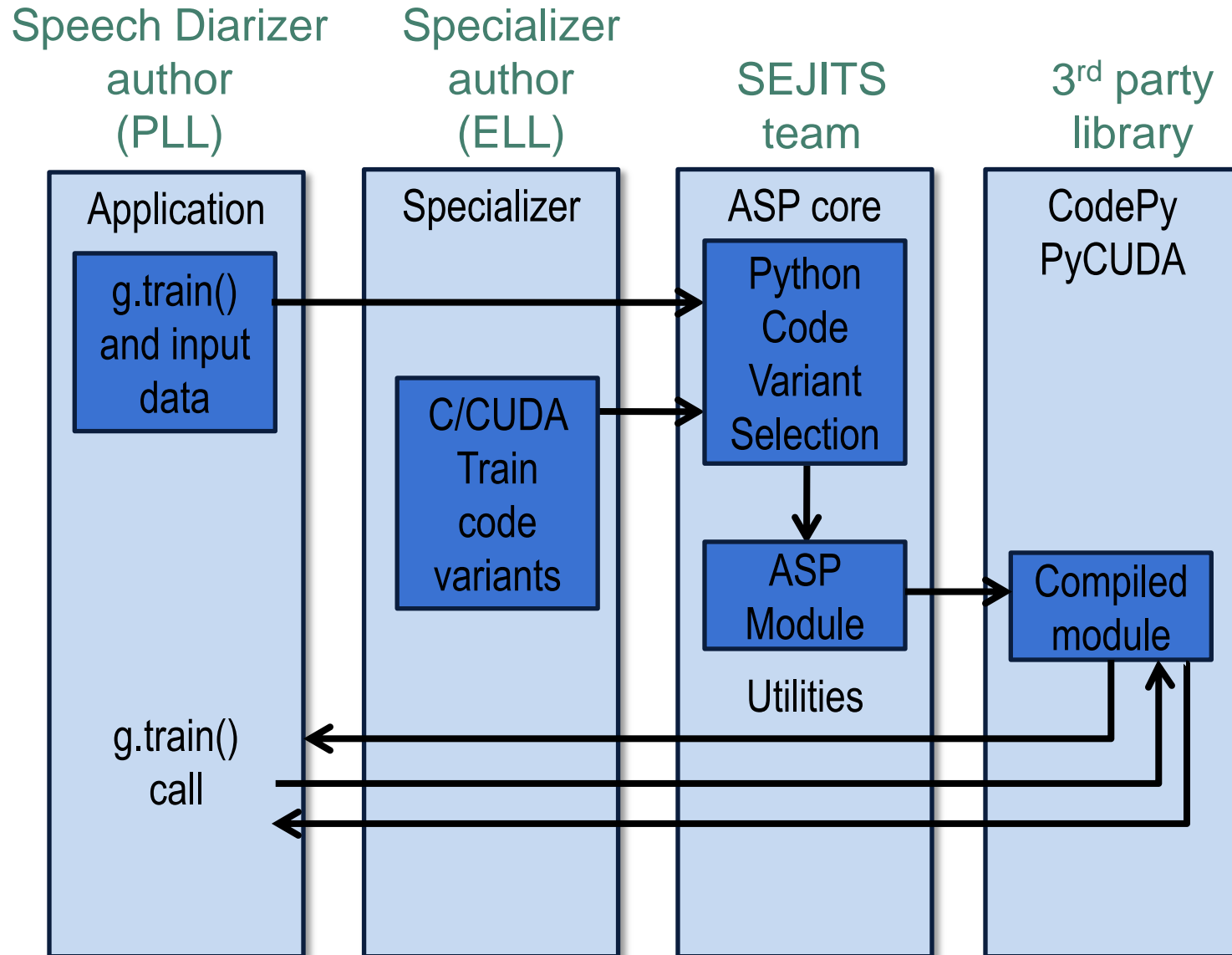
SEJITS Framework

- Python code that handles application
 - Manipulates problem data, determines learning targets
- C/CUDA code that runs quickly
 - Allocates GPU memory
 - Performs main EM iterative loop
- Specializer (ASP)
 - Selects appropriate code variant (from history)
 - Pulls in the template for the code variant, parameterizes it and compiles to binary





Separation of Concerns





C

```
# Get the per_cluster GMM on a cluster
init_train_gmm(self, N, per_cluster)))

new_gmm_list(M,D)

# Perform hierarchical agglomeration based on BIC scores
best_BIC_score = 1.0
while (best_BIC_score > 0 and len(self.gmm_list) > 1):

    num_clusters = len(self.gmm_list)

    # Resegment data based on likelihood scoring
    likelihoods = self.gmm_list[0].score(self.X)
    for g in self.gmm_list[1:]:
        likelihoods = np.column_stack((likelihoods, g.score(self.X)))
    most_likely = likelihoods.argmax(axis=1)

    # Across 2.5 secs of observations, vote on which cluster they should be associated with
    split_votes = self.gmm_list[0].split_votes(most_likely, self.X)

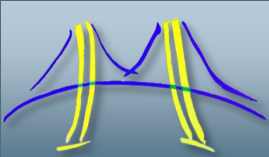
    for g in self.gmm_list[1:]:
        g.train(x)

    # Score all pairs of GMMs using BIC
    best_merged_gmm = None
    best_BIC_score = 0.0
    merged_tuple = None

    for gmm1idx in range(len(Iter_bic_list)):
        for gmm2idx in range(gmm1idx+1, len(Iter_bic_list)):
            g1, d1 = Iter_bic_list[gmm1idx]
            g2, d2 = Iter_bic_list[gmm2idx]
            score = 0.0
            new_gmm, score = compute_distance_BIC(g1, g2, np.concatenate((d1, d2)))
            if score > best_BIC_score:
                best_merged_gmm = new_gmm
                merged_tuple = (g1, g2)
                best_BIC_score = score

    # Merge the winning candidate pair
    if best_BIC_score > 0.0:
        self.gmm_list.remove(merged_tuple[0])
        self.gmm_list.remove(merged_tuple[1])
        self.gmm_list.append(best_merged_gmm)
```

[illegible]



Speaker Diarization in Python

Python

```
def AHC(self):
```

```
# Get the events, divide them into an initial k clusters and train each GMM on a cluster
per_cluster = self.N/self.init_num_clusters
init_training = zip(self.gmm_list,np.vsplit(self.X, range(per_cluster, self.N, per_cluster)))
for g, x in init_training:
    g.train(x)
```

```
# Perform hierarchical agglomeration based on BIC scores
best_BIC_score = 1.0
while (best_BIC_score > 0 and len(self.gmm_list) > 1):
```

```
    num_clusters = len(self.gmm_list)
```

```
# Resegment data based on likelihood scoring
likelihoods = self.gmm_list[0].score(self.X)
for g in self.gmm_list[1:]:
    likelihoods = np.column_stack((likelihoods, g.score(self.X)))
most_likely = likelihoods.argmax(axis=1)
```

```
# Across 2.5 secs of observations, vote on which
split_events = split_events_based_on_votes(most_likely,
```

```
for g, data in split_events:
    g.train(data)
```

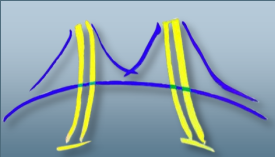
```
# Score all pairs of GMMs using BIC
best_merged_gmm = None
best_BIC_score = 0.0
merged_tuple = None
```

```
for gmmidx in range(len(iter_bic_list)):
    for gmm2idx in range(gmmidx+1, len(iter_bic_list)):
        g1, d1 = iter_bic_list[gmmidx]
        g2, d2 = iter_bic_list[gmm2idx]
        score = 0.0
        new_gmm, score = compute_distance_BIC(g1, g2, np.concatenate((d1, d2)))
        if score > best_BIC_score:
            best_merged_gmm = new_gmm
            merged_tuple = (g1, g2)
            best_BIC_score = score
```

```
# Merge the winning candidate pair
if best_BIC_score > 0.0:
    self.gmm_list.remove(merged_tuple[0])
    self.gmm_list.remove(merged_tuple[1])
    self.gmm_list.append(best_merged_gmm)
```

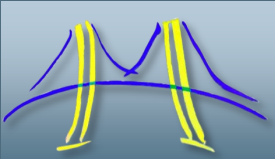
C

```
/*
 * Copyright (c) 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108, 2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121, 2122, 2123, 2124, 2125, 2126, 2127, 2128, 2129, 2130, 2131, 2132, 2133, 2134, 2135, 2136, 2137, 2138, 2139, 2140, 2141, 2142, 2143, 2144, 2145, 2146, 2147, 2148, 2149, 2150, 2151, 2152, 2153, 2154, 2155, 2156, 2157, 2158, 2159, 2160, 2161, 2162, 2163, 2164, 2165, 2166, 2167, 2168, 2169, 2170, 2171, 2172, 2173, 2174, 2175, 2176, 2177, 2178, 2179, 2180, 2181, 2182, 2183, 2184, 2185, 2186, 2187, 2188, 2189, 2190, 2191, 2192, 2193, 2194, 2195, 2196, 2197, 2198, 2199, 2200, 2201, 2202, 2203, 2204, 2205, 2206, 2207, 2208, 2209, 2210, 2211, 2212, 2213, 2214, 2215, 2216, 2217, 2218, 2219, 2220, 2221, 2222, 2223, 2224, 2225, 2226, 2227, 2228, 2229, 2230, 2231, 2232, 2233, 2234, 2235, 2236, 2237, 2238, 2239, 2240, 2241, 2242, 2243, 2244, 2245, 2246, 2247, 2248, 2249, 2250, 2251, 2252, 2253, 2254, 2255, 2256, 2257, 2258, 2259, 2260, 2261, 2262, 2263, 2264, 2265, 2266, 2267, 2268, 2269, 2270, 2271, 2272, 2273, 2274, 2275, 2276, 2277, 2278, 2279, 2280, 2281, 2282, 2283, 2284, 2285, 2286, 2287, 2288, 2289, 2290, 2291, 2292, 2293, 2294, 2295, 2296, 2297, 2298, 2299, 2300, 2301, 2302, 2303, 2304, 2305, 2306, 2307, 2308, 2309, 2310, 2311, 2312, 2313, 2314, 2315, 2316, 2317, 2318, 2319, 2320, 2321, 2322, 2323, 2324, 2325, 2326, 2327, 2328, 2329, 2330, 2331, 2332, 2333, 2334, 2335, 2336, 2337, 2338, 2339, 2340, 2341, 2342, 2343, 2344, 2345, 2346, 2347, 2348, 2349, 2350, 2351, 2352, 2353, 2354, 2355, 2356, 2357, 2358, 2359, 2360, 2361, 2362, 2363, 2364, 2365, 2366, 2367, 2368, 2369, 2370, 2371, 2372, 2373, 2374, 2375, 2376, 2377, 2378, 2379, 2380, 2381, 2382, 2383, 2384, 2385, 2386, 2387, 2388, 2389, 2390, 2391, 2392, 2393, 2394, 2395, 2396, 2397, 2398, 2399, 2400, 2401, 2402, 2403, 2404, 2405, 2406, 2407, 2408, 2409, 2410, 2411, 2412, 2413, 2414, 2415, 2416, 2417, 2418, 2419, 2420, 2421, 2422, 2423, 2424, 2425, 2426, 2427, 2428, 2429, 2430, 2431, 2432, 2433, 2434, 2435, 2436, 2437, 2438, 2439, 2440, 2441, 2442, 2443, 2444, 2445, 2446, 2447, 2448, 2449, 2450, 2451, 2452, 2453, 2454, 2455, 2456, 2457, 2458, 2459, 2460, 2461, 2462, 2463, 2464, 2465, 2466, 2467, 2468, 2469, 2470, 2471, 2472, 2473, 2474, 2475, 2476, 2477, 2478, 2479, 2480, 2481, 2482, 2483, 2484, 2485, 2486, 2487, 2488, 2489, 2490, 2491, 2492, 2493, 2494, 2495, 2496, 2497, 2498, 2499, 2500, 2501, 2502, 2503, 2504, 2505, 2506, 2507, 2508, 2509, 2510, 2511, 2512, 2513, 2514, 2515, 2516, 2517, 2518, 2519, 2520, 2521, 2522, 2523, 2524, 2525, 2526, 2527, 2528, 2529, 2530, 2531, 2532, 2533, 2534, 2535, 2536, 2537, 2538, 2539, 2540, 2541, 2542, 2543, 2544, 2545, 2546, 2547, 2548, 2549, 2550, 2551, 2552, 2553, 2554, 2555, 2556, 2557, 2558, 2559, 2560, 2561, 2562, 2563, 2564, 2565, 2566, 2567, 2568, 2569, 2570, 2571, 2572, 2573, 2574, 2575, 2576, 2577, 2578, 2579, 2580, 2581, 2582, 2583, 2584, 2585, 2586, 2587, 2588, 2589, 2590, 2591, 2592, 2593, 2594, 2595, 2596, 2597, 2598, 2599, 2600, 2601, 2602, 2603, 2604, 2605, 2606, 2607, 2608, 2609, 2610, 2611, 2612, 2613, 2614, 2615, 2616, 2617, 2618, 2619, 2620, 2621, 2622, 2623, 2624, 2625, 2626, 2627, 2628, 2629, 2630, 2631, 2632, 2633, 2634, 2635, 2636, 2637, 2638, 2639, 2640, 2641, 2642, 2643, 2644, 2645, 2646, 2647, 2648, 2649, 2650, 2651, 2652, 2653, 2654, 2655, 2656, 2657, 2658, 2659, 2660, 2661, 2662, 2663, 2664, 2665, 2666, 2667, 2668, 2669, 2670, 2671, 2672, 2673, 2674, 2675, 2676, 2677, 2678, 2679, 2680, 2681, 2682, 2683, 2684, 2685, 2686, 2687, 2688, 2689, 2690, 2691, 2692, 2693, 2694, 2695, 2696, 2697, 2698, 2699, 2700, 2701, 2702, 2703, 2704, 2705, 2706, 2707, 2708, 2709, 2710, 2711, 2712, 2713, 2714, 2715, 2716, 2717, 2718, 2719, 2720, 2721, 2722, 2723, 2724, 2725, 2726, 2727, 2728, 2729, 2730, 2731, 2732, 2733, 2734, 2735, 2736, 2737, 2738, 2739, 2740, 2741, 2742, 2743, 2744, 2745, 2746, 2747, 2748, 2749, 2750, 2751, 2752, 2753, 2754, 2755, 2756, 2757, 2758, 2759, 2760, 2761, 2762, 2763, 2764, 2765, 2766, 2767, 2768, 2769, 2770, 2771, 2772, 2773, 2774, 2775, 2776, 2777, 2778, 2779, 2780, 2781, 2782, 2783, 2784, 2785, 2786, 2787, 2788, 2789, 2790, 2791, 2792, 2793, 2794, 2795, 2796, 2797, 2798, 2799, 2800, 2801, 2802, 2803, 2804, 2805, 2806, 2807, 2808, 2809, 2810, 2811, 2812, 2813, 2814, 2815, 2816, 2817, 2818, 2819, 2820, 2821, 2822, 2823, 2824, 2825, 2826, 2827, 2828, 2829, 2830, 2831, 2832, 2833, 2834, 2835, 2836, 2837, 2838, 2839, 2840, 2841, 2842, 2843, 2844, 2845, 2846, 2847, 2848, 2849, 2850, 2851, 2852, 2853, 2854, 2855, 2856, 2857, 2858, 2859, 2860, 2861, 2862, 2863, 2864, 2865, 2866, 2867, 2868, 2869, 2870, 2871, 2872, 2873, 2874, 2875, 2876, 2877, 2878, 2879, 2880, 2881, 2882, 2883, 2884, 2885, 2886, 2887, 2888, 2889, 2890, 2891, 2892, 2893, 2894, 2895, 2896, 2897, 2898, 2899, 2900, 2901, 2902, 2903, 2904, 2905, 2906, 2907, 2908, 2909, 2910, 2911, 2912, 2913, 2914, 2915, 2916, 2917, 2918, 2919, 2920, 2921, 2922, 2923, 2924, 2925, 2926, 2927, 2928, 2929, 2930, 2931, 2932, 2933, 2934, 2935, 2936, 2937, 2938, 2939, 2940, 2941, 2942, 2943, 2944, 2945, 2946, 2947, 2948, 2949, 2950, 2951, 2952, 2953, 2954, 2955, 2956, 2957, 2958, 2959, 2960, 2961, 2962, 2963, 2964, 2965, 2966, 2967, 2968, 2969, 2970, 2971, 2972, 2973, 2974, 2975, 2976, 2977, 2978, 2979, 2980, 2981, 2982, 2983, 2984, 2985, 2986, 2987, 2988, 2989, 2990, 2991, 2992, 2993, 2994, 2995, 2996, 2997, 2998, 2999, 3000, 3001, 3002, 3003, 3004, 3005, 3006, 3007, 3008, 3009, 3010, 3011, 3012, 3013, 3014, 3015, 3016, 3017, 3018, 3019, 3020, 3021, 3022, 3023, 3024, 3025, 3026, 3027, 3028, 3029, 3030, 3031, 3032, 3033, 3034, 3035, 3036, 3037, 3038, 3039, 3040, 3041, 3042, 3043, 3044, 3045, 3046, 3047, 3048, 3049, 3050, 3051, 3052, 3053, 3054, 3055, 3056, 3057, 3058, 3059, 3060, 3061, 3062, 3063, 3064, 3065, 3066, 3067, 3068, 3069, 3070, 3071, 3072, 3073, 3074, 3075, 3076, 3077, 3078, 3079, 3080, 3081, 3082, 3083, 3084, 3085, 3086, 3087, 3088, 3089, 3090, 3091, 3092, 3093, 3094, 3095, 3096, 3097, 3098, 3099, 3100, 3101, 3102, 3103, 3104, 3105, 3106, 3107, 3108, 3109, 3110, 3111, 3112, 3113, 3114, 3115, 3116, 3117, 3118, 3119, 3120, 3121, 3122, 3123, 3124, 3125, 3126, 3127, 3128, 3129, 3130, 3131, 3132, 3133, 3134, 3135, 3136, 3137, 3138, 3139, 3140, 3141, 3142, 3143, 3144, 3145, 3146, 3147, 3148, 3149, 3150, 3151, 3152, 3153, 3154, 3155, 3156, 3157, 3158, 3159, 3160, 3161, 3162, 3163, 3164, 3165, 3166, 3167, 3168, 3169, 3170, 3171, 3172, 3173, 3174, 3175, 3176, 3177, 3178, 3179, 3180, 3181, 3182, 3183, 3184, 3185, 3186, 3187, 3188, 3189, 3190, 3191, 3192, 3193, 3194, 3195, 3196, 3197, 3198, 3199, 3200, 3201, 3202, 3203, 3204, 3205, 3206, 3207, 3208, 3209, 3210, 3211, 3212, 3213, 3214, 3215, 3216, 3217, 3218, 3219, 3220, 3221, 3222, 3223, 3224, 3225, 3226, 3227, 3228, 3229, 3230, 3231, 3232, 3233, 3234, 3235, 3236, 3237, 3238, 3239, 3240, 3241, 3242, 3243, 3244, 3245, 3246, 3247, 3248, 3249, 3250, 3251, 3252, 3253, 3254, 3255, 3256, 3257, 3258, 3259, 3260, 3261, 3262, 3263, 3264, 3265, 3266, 3267, 3268, 3269, 3270, 3271, 3272, 3273, 3274, 3275, 3276, 3277, 3278, 3279, 3280, 3281, 3282, 3283, 3284, 3285, 3286, 3287, 3288, 3289, 3290, 3291, 3292, 3293, 3294, 3295, 3296, 3297, 3298, 3299, 3300, 3301, 3302, 3303, 3304, 3305, 3306, 3307, 3308, 3309, 3310, 3311, 3312, 3313, 3314, 3315, 3316, 3317, 3318, 3319, 3320, 3321, 3322, 3323, 3324, 3325, 3326, 3327, 3328, 3329, 3330, 3331, 3332, 3333, 3334, 3335, 3336, 3337, 3338, 3339, 3340, 3341, 3342, 3343, 3344, 3345, 3346, 3347, 3348, 3349, 3350, 3351, 3352, 3353, 3354, 3355, 3356, 3357, 3358, 3359, 3360, 3361, 3362, 3363, 3364, 3365, 3366, 3367, 3368, 3369, 3370, 3371, 3372, 3373, 3374, 3375, 3376, 3377, 3378, 3379, 3380, 3381, 3382, 3383, 3384, 3385, 3386, 3387, 3388, 3389, 3390, 3391, 3392, 3393, 3394, 3395, 3396, 3397, 3398, 3399, 3400, 3401, 3402, 3403, 3404, 3405, 3406, 3407, 3408, 3409, 3410, 3411, 3412, 3413, 3414, 3415, 3416, 3417, 3418, 3419, 3420, 3421, 3422, 3423, 3424, 3425, 3426, 3427, 3428, 3429, 3430, 3431, 3432, 3433, 3434, 3435, 3436, 3437, 3438, 3439, 3440, 3441, 3442, 3443, 3444, 3445, 3446, 3447, 3448, 3449, 3450, 3451, 3452, 3453, 3454, 3455, 3456, 3457, 3458, 3459, 3460, 3461, 3462, 3463, 3464, 3465, 3466, 3467, 3468, 3469, 3470, 3471, 3472, 3473, 3474, 3475, 3476, 3477, 3478, 3479, 3480, 3481, 3482, 3483, 3484, 3485, 3486, 3487, 3488, 3489, 3490, 3491, 3492, 3493, 3494, 3495, 3496, 3497, 3498, 3499, 3500, 3501, 3502, 3503, 3504, 3505, 3506, 3507, 3508, 3509, 3510, 3511, 3512, 3513, 3514, 3515, 3516, 3517, 3518, 3519, 3520, 3521, 3522, 3523, 3524, 3525, 3526, 3527, 3528, 3529, 3530, 3531, 3532, 3533, 3534, 3535, 3536, 3537, 3538, 3539, 3540, 3541, 3542, 3543, 3544, 3545, 3546, 3547, 3548, 3549, 3550, 3551, 3552, 3553, 3554, 3555, 3556, 3557, 3558, 3559, 3560, 3561, 3562, 3563, 3564, 3565, 3566, 3567, 3568, 3569, 3570, 3571, 3572, 3573, 3574, 3575, 3576, 3577, 3578, 3579, 3580, 3581, 3582, 3583, 3584, 3585, 3586, 3587, 3588, 3589, 3590, 3591, 3592, 3593, 3594, 3595, 3596, 3597, 3598, 3599, 3600, 3601, 3602, 3603, 3604, 3605, 3606, 3607, 3608, 3609, 3610, 3611, 3612, 3613, 3614, 3615, 3616, 3617, 3618, 3619, 3620, 3621, 3622, 3623, 3624, 3625, 3626, 3627, 3628, 3629, 3630, 3631, 3632, 3633, 3634, 3635, 3636, 3637, 3638, 3639, 3640, 3641, 3642, 3643, 3644, 3645, 3646, 3647, 3648, 3649, 3650, 3651, 3652, 3653, 3654, 3655, 3656, 3657, 3658, 3659, 3660, 3661, 3662, 3663, 3664, 3665, 3666, 3667, 3668, 3669, 3670, 3671, 3672, 3673, 3674, 3675, 3676, 3677, 3678, 3679, 3680, 3681, 3682, 3683, 3684, 3685, 3686, 3687, 3688, 3689, 3690, 3691, 3692, 3693, 3694, 3695, 3696, 3697, 3698, 3699, 3700, 3701, 3702, 3703, 3704, 3705, 3706, 3707, 3708, 3709, 3710, 3711, 3712, 3713, 3714, 3715, 3716, 3717, 3718, 3719, 3720, 3721, 3722, 3723, 3724, 3725, 3726, 3727, 3728, 3729, 3730, 3731, 3732, 3733, 3734, 3735, 3736, 3737, 3738, 3739, 3740, 3741, 3742, 3743, 3744, 3745, 3746, 3747, 3748, 3749, 3750, 3751, 3752, 3753, 3754, 3755, 3756, 3757, 3758, 3759, 3760, 3761, 3762, 3763, 3764, 3765, 3766, 3767, 3768, 3769, 3770, 3771, 3772, 3773, 3774, 3775, 3776, 3777, 3778, 3779, 3780, 3781, 3782, 3783, 3784, 3785, 3786, 3787, 3788, 3789, 3790, 3791, 3792, 3793, 3794, 3795, 3796, 3797, 3798, 3799, 3800, 3801, 3802, 3803, 3804, 3805, 3806, 3807, 3808, 3809, 3810, 3811, 3812, 3813, 3814, 3815, 3816, 3817, 3818, 3819, 3820, 3821, 3822, 3823, 3824, 3825, 3826, 3827, 3828, 3829, 3830, 3831, 3832, 3833, 3834, 3835, 3836, 3837, 3838, 3839, 3840, 3841, 3842, 3843, 3844, 3845, 3846, 3847, 3848, 3849, 3850, 3851, 3852, 3853, 3854, 3855, 3856, 3857, 3858, 3859, 3860, 3861, 3862, 3863, 3864, 3865, 3866, 3867, 3868, 3869, 3870, 3871, 3872, 3873, 3874, 3875, 3876, 3877, 3878, 3879, 3880, 3881, 3882, 3883, 3884, 3885, 3886, 3887, 3888, 3889, 3890, 3891, 3892, 3893, 3894, 3895, 3896, 3897, 3898, 3899, 3900, 3901, 3902, 3903, 3904, 3905, 3906, 3907, 3908, 3909, 3910, 3911, 3912,
```



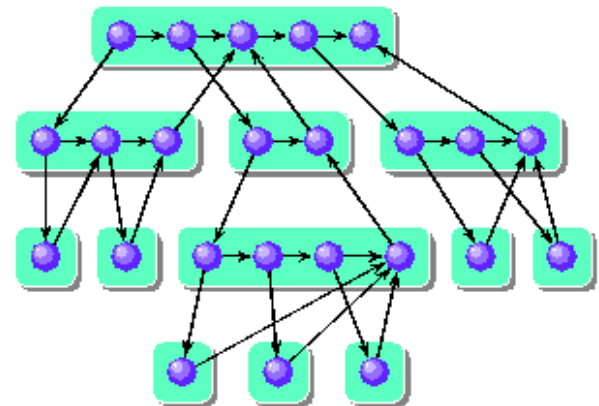
Results – Specializer Overhead

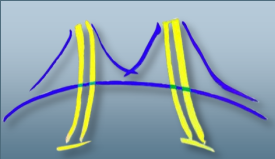
- Python AHC code is within **1.25x** of pure C/CUDA implementation performance
- C/CUDA AHC (from winter retreat) – 250x realtime
- SEJITized AHC ~ 200x realtime
- Time lost in:
 - Outer loop and GMM creation in Python
 - Data copying overhead from CPU to GPU
 - GMM scoring in Python



Cilk Backend

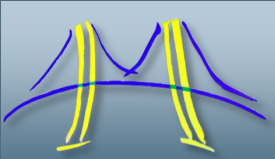
- We have implemented the Cilk backend for GMM training
- ASP selects version based on available hardware
- Current implementation ~100x realtime
- 5-10% C code reused
- All specializer infrastructure reused





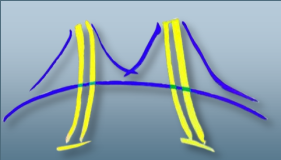
Results – Portability & Maintenance

- Specializes to two types of platforms (multi-core CPU, Nvidia GPU) to support **portability**
 - Exact same application code
- **Reuse** of infrastructure:
 - Specializer creation and code variant selection mechanism reused
- **Maintaining the code** for next generation of hardware
 - Task of specializer writer, transparent to the application developer



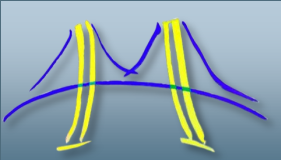
Conclusion & Future Work

- SEJITized GMM training in Speaker Diarization component of Meeting Diarist
- Specialized covariance matrix computation with code variant selection to two platforms
- Currently a factor of 1.25x slower than pure C/CUDA implementation (200 x faster than realtime)
- Future work:
 - Further specialize train kernel
 - SEJITize other components
 - Improve code variant selection mechanism

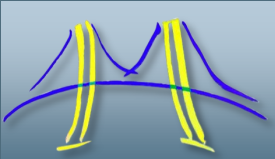


Thank you!

Questions?

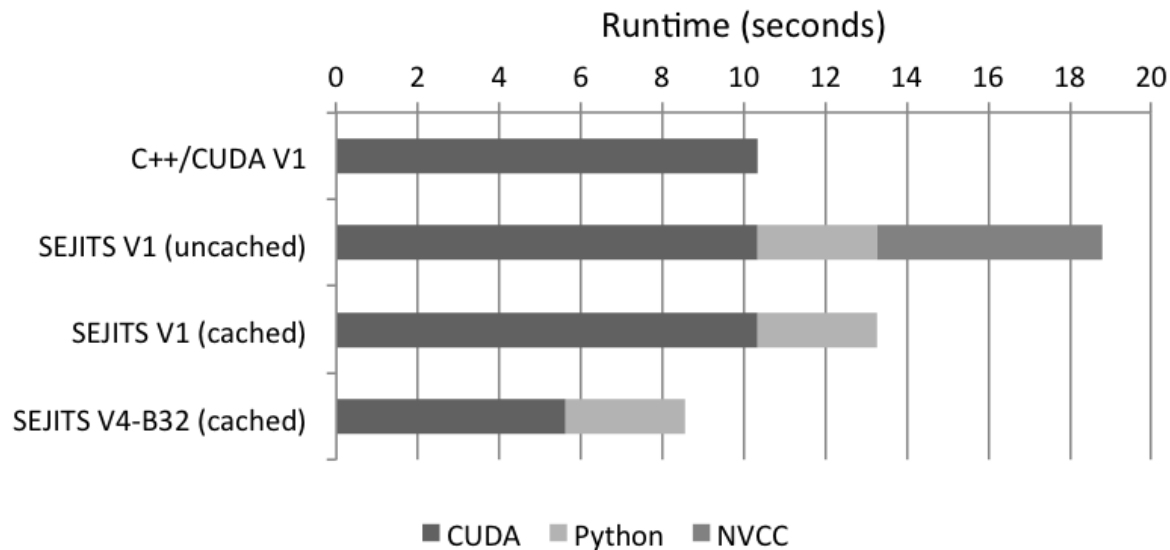


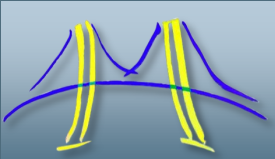
Backup Slides



Results – Specializer Overhead in AHC

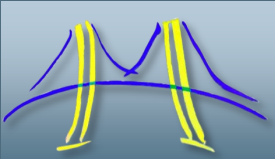
- Initial invocation – 81% overhead due to compiler invocations
- Future runs using automatically determined optimal code variant achieve 17% performance *improvement* over the original GPU implementation (V1)





ASP vs Auto-tuning Libraries

	ATLAS	FFTW, Spiral, OSKI	ASP/GM M	ASP/Sten cil	Delite/ OptiML	Copperhe ad
Autuning of code	✓	✓	✓	✓	✓	✓
Based on runtime information		✓	✓	✓	✓	✓
Based on higher- order func				✓	✓	✓
Using reusable framework			✓	✓	✓	✓
Embedded in HLL			✓	✓	✓	✓

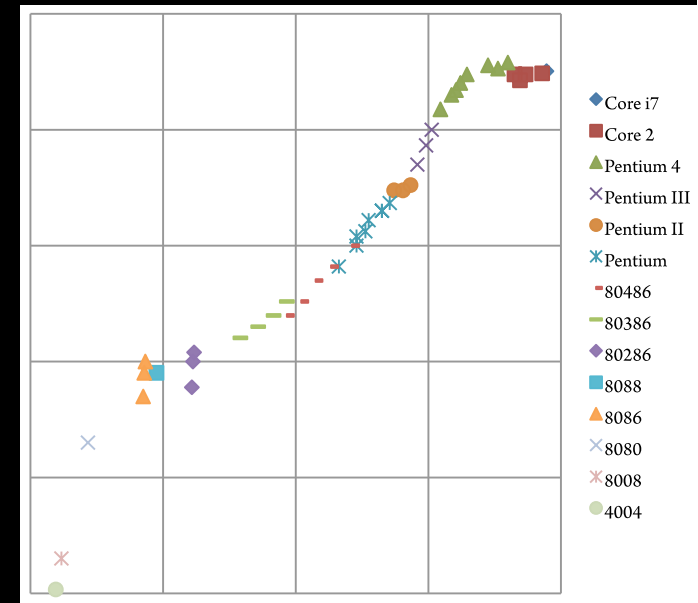


The shift to parallel processing

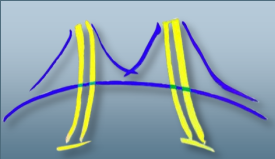
■ Parallel processing is here

“This shift toward increasing parallelism is not a triumphant stride forward based on breakthroughs in novel software and architectures for parallelism; instead, this plunge into parallelism is actually a retreat from even greater challenges that thwart efficient silicon implementation of traditional uniprocessor architectures.”

- The Berkeley View

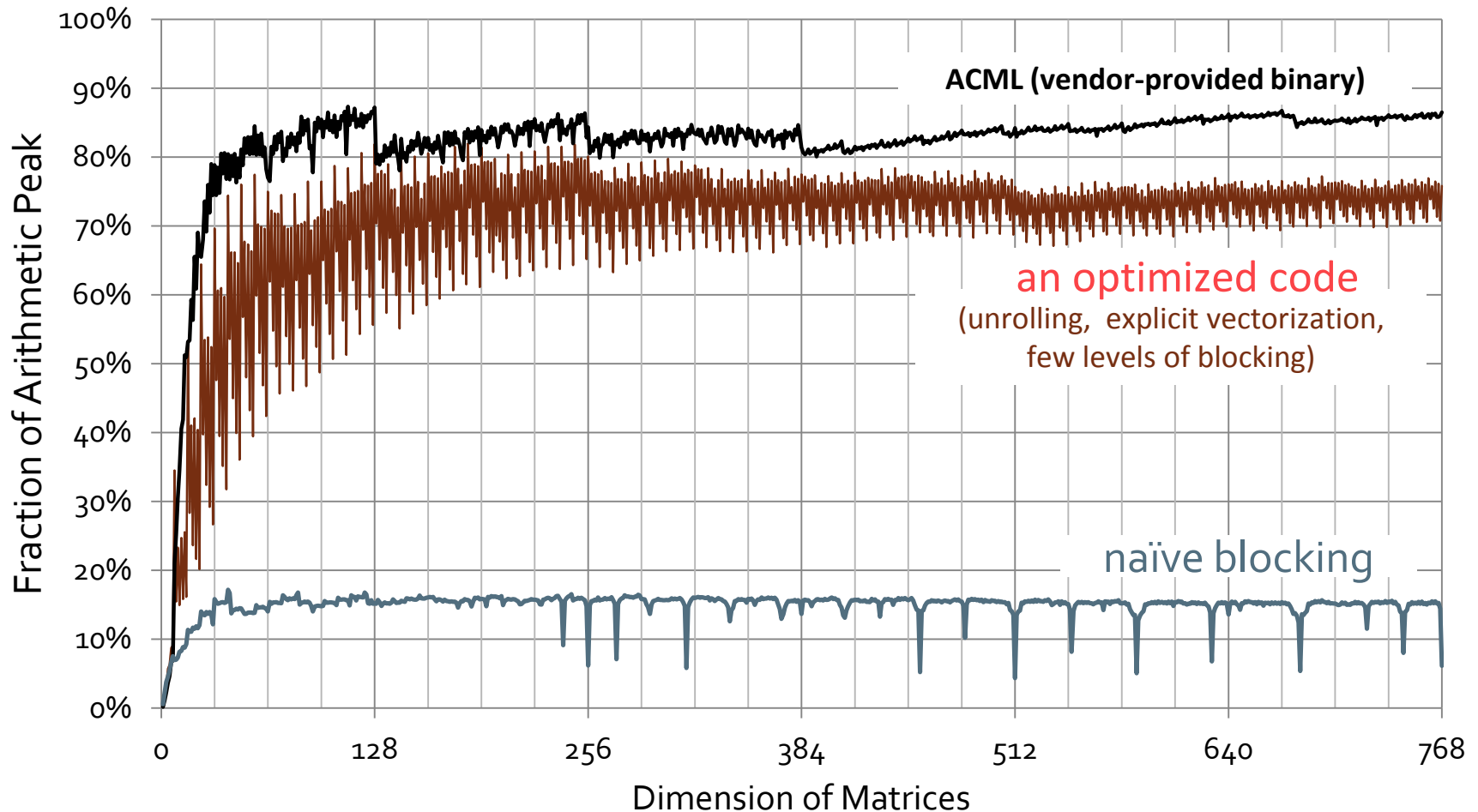


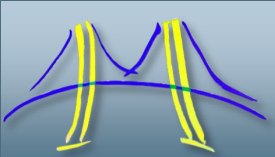
Intel Processor Clock Speed



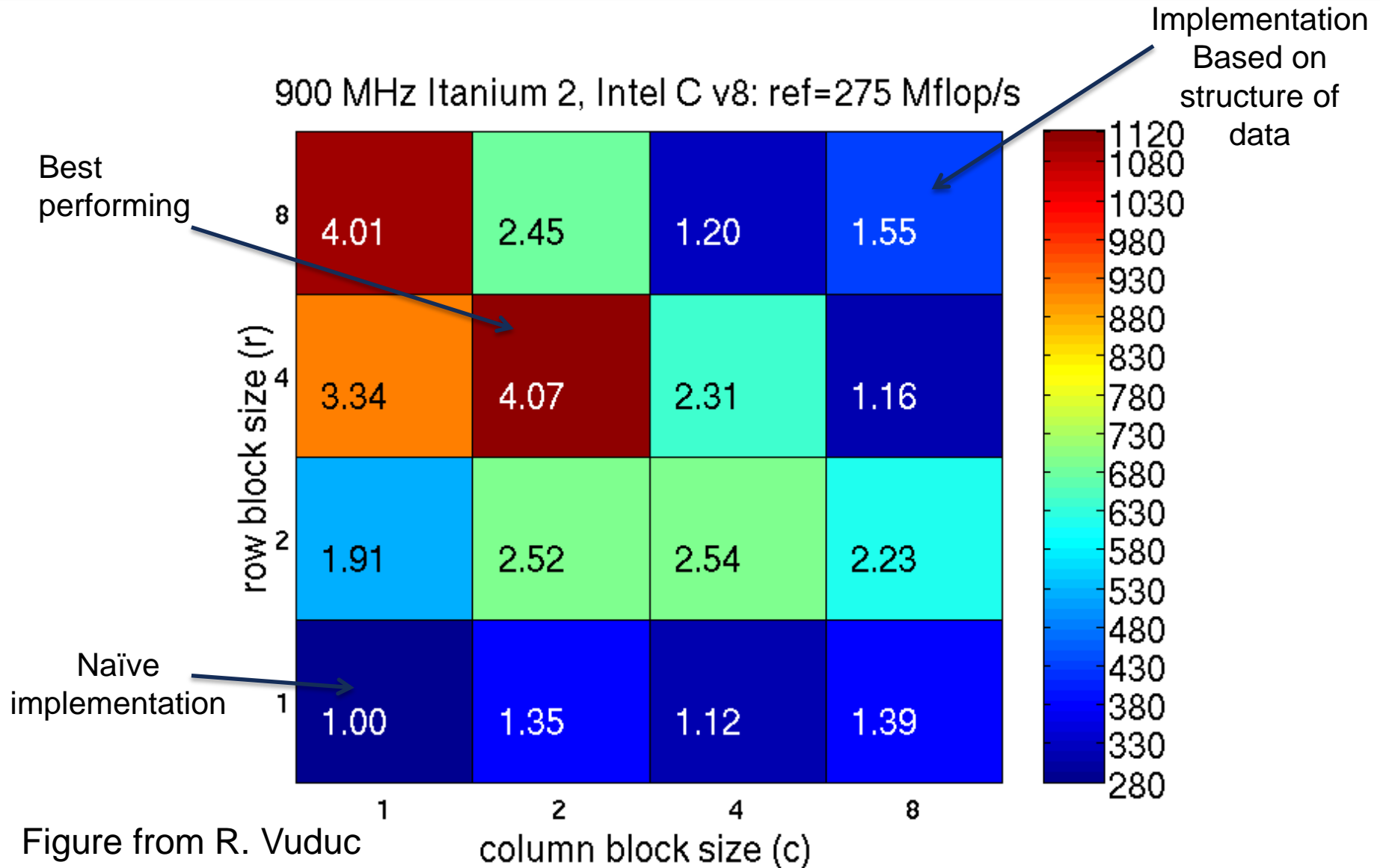
Writing Fast Code is Hard

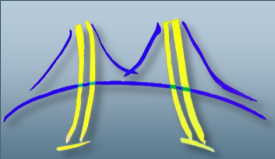
Dense Matrix Multiply (V. Volkov)





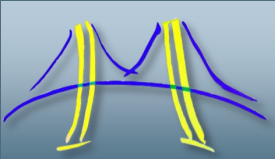
Finding Best Implementation is Hard



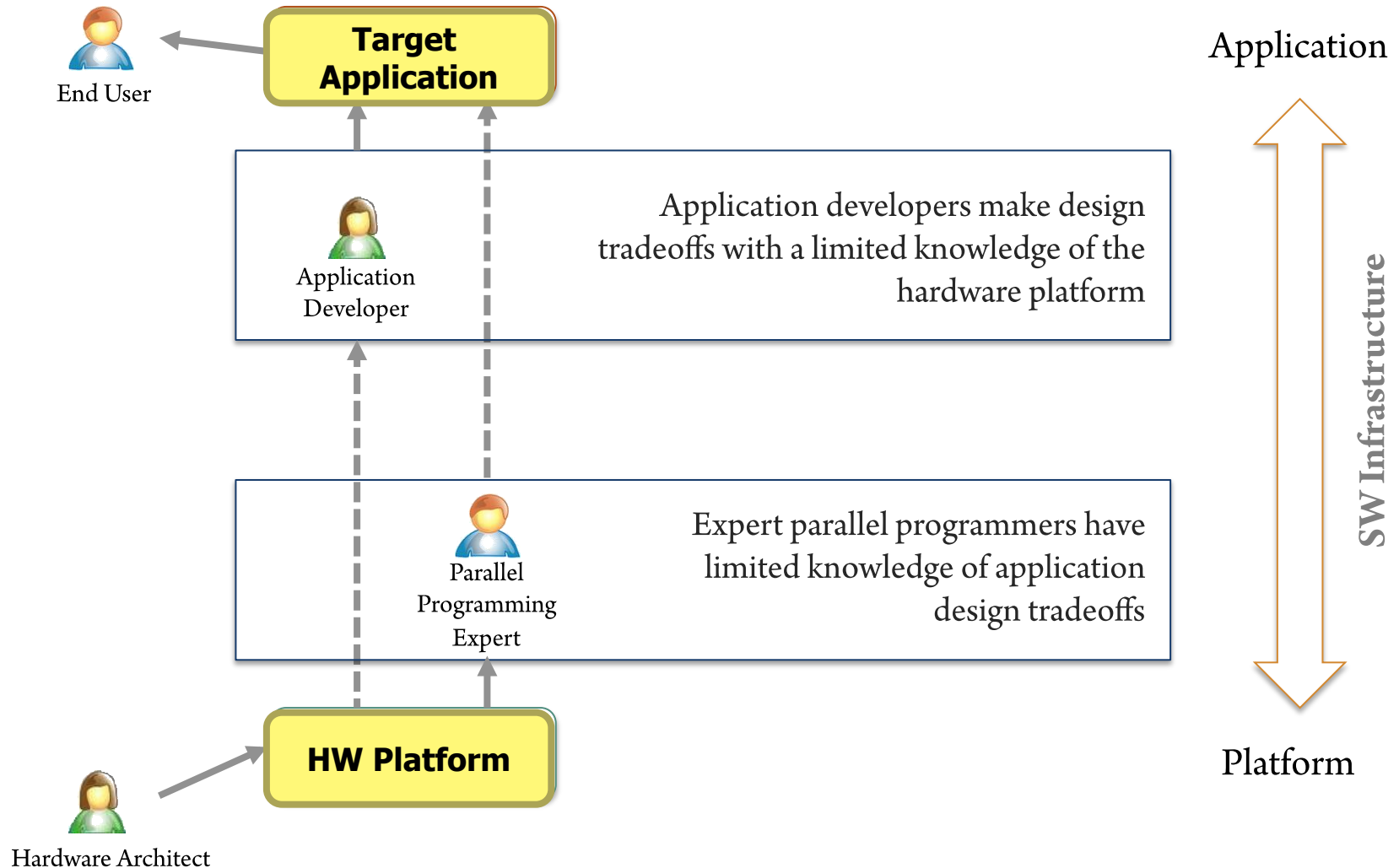


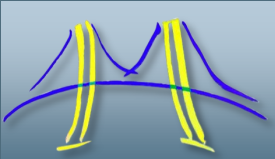
Productivity vs Performance

- Scientists and domain experts prefer to use high-level languages such as Python or MATLAB
- However, to achieve sufficient performance, computationally-intensive parts of applications must eventually be rewritten in low-level languages
- In addition, parallel platform details and input parameters determine the best-performing parallel implementation



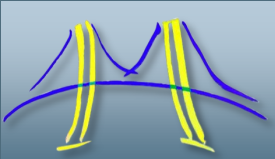
Implementation Gap





Outline

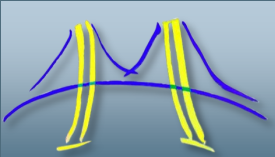
- SEJITS approach
- Gaussian Mixture Model & Applications
- Covariance Matrix Computation & Code Variants
- Specialization
- Results
- Conclusion & Future Work



Selective Embedded Just-In-Time Specialization (SEJITS)

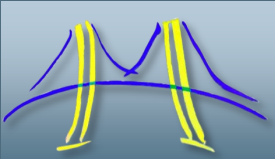
Key Idea: Generate, compile, and execute high performance parallel code at runtime using code transformation, introspection, variant selection and other features of high-level languages.

Invisibly to the user.

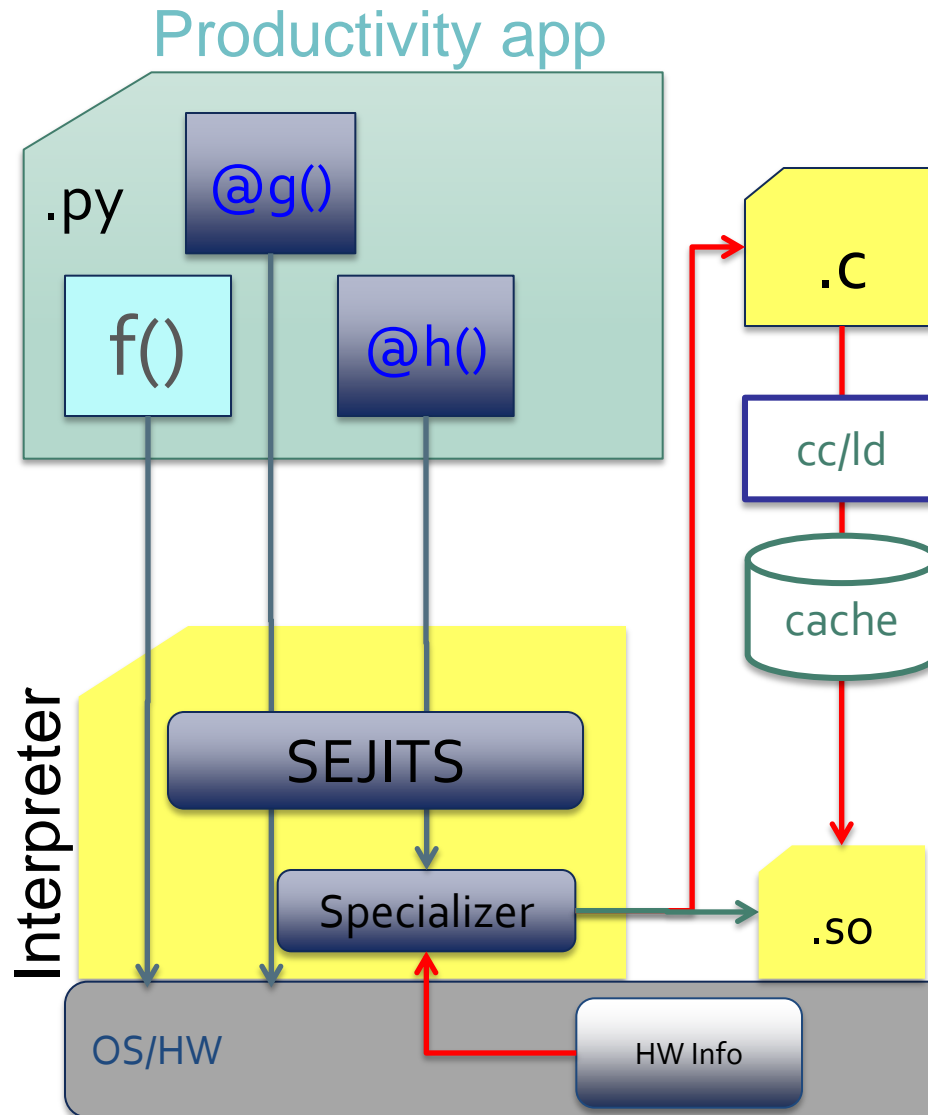


Selective Embedded JIT Specialization (SEJITS)

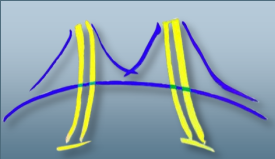
- Leverage patterns to bridge productivity and efficiency
- PLL (productivity-level language, eg Python) for applications
- “Specializers” generate ELL (efficiency-level language) code targeted to hardware
 - Code generation can happen at runtime
 - Specializers can incorporate autotuning
 - Think: pattern-specific embedded DSLs
- ELL performance with PLL effort



Selective Embedded JIT Specialization (SEJITS)



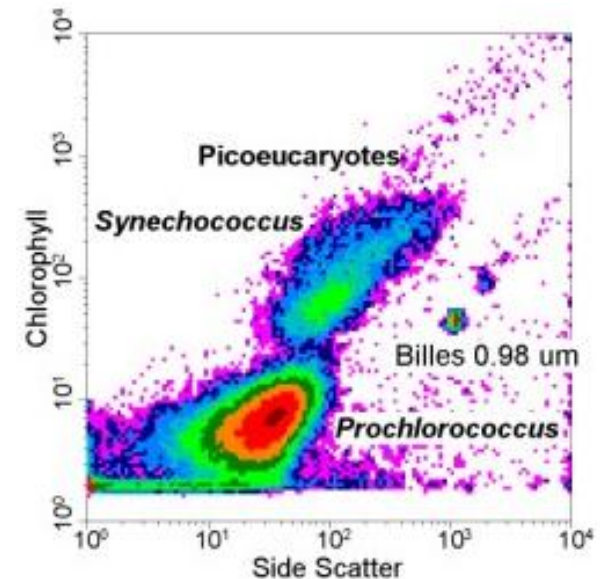
ASP – A SEJITS
for Python

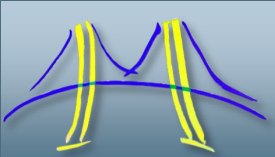


Applications of Gaussian Mixture Models

► Applications

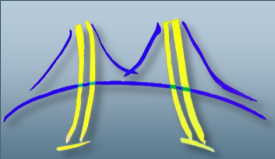
- Can be used to cluster/classify any sequence of observations
- Speech Recognition – speaker classification, acoustic modeling for speech recognition
- Computer Vision – image segmentation, hand writing recognition
- Biology – flow cytometry
- Data mining – topic classification in web documents
- Many more...





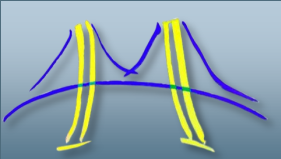
Results – Specializer Overhead

- Application example – Agglomerative Hierarchical Clustering for Speaker Diarization
 - Uses GMMs to represent distribution of audio features for speakers in a recorded meeting
 - Iteratively trains GMMs using different number of components each time and measuring which number of components best fits the data
 - Number of components in the best GMM corresponds to number of speakers in the meeting



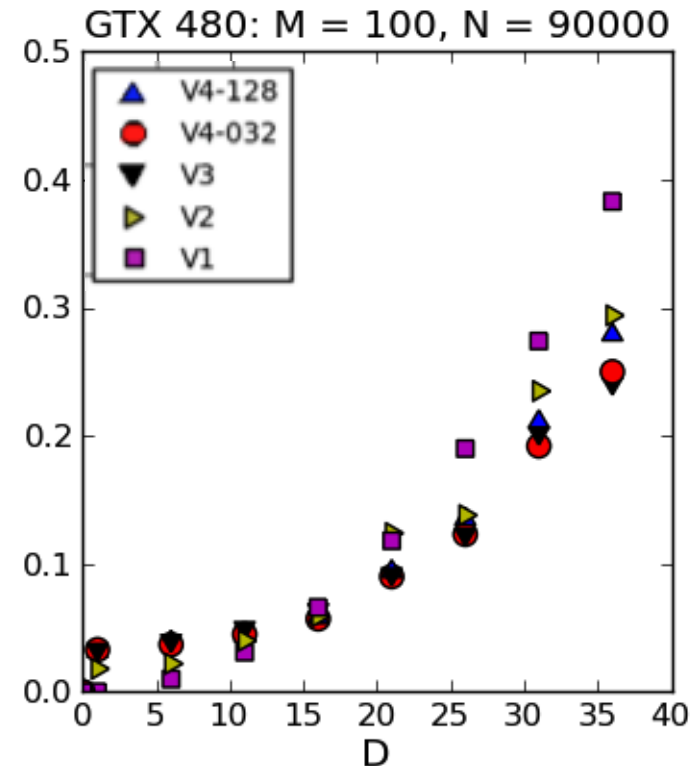
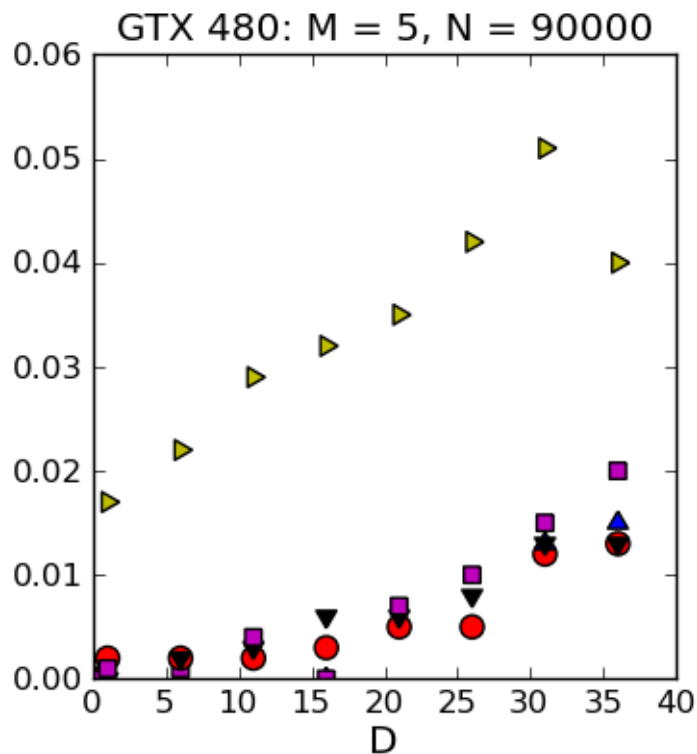
Conclusions & Future Work

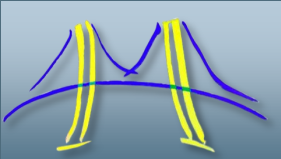
- ASP framework encapsulates code variant selection mechanisms and handcrafted templates to:
 - Allow domain expert to stay in the high-level language domain and focus on the application
 - Obtain high performance from expert-tuned code
- Example in Gaussian Mixture Model Applications
- Performance benefit of specialization outweighs the overhead of Python and the JIT process
- Expand to:
 - more platforms, applications, patterns
 - other code variant selection mechanisms



Results – Version Comparison (Raw CUDA)

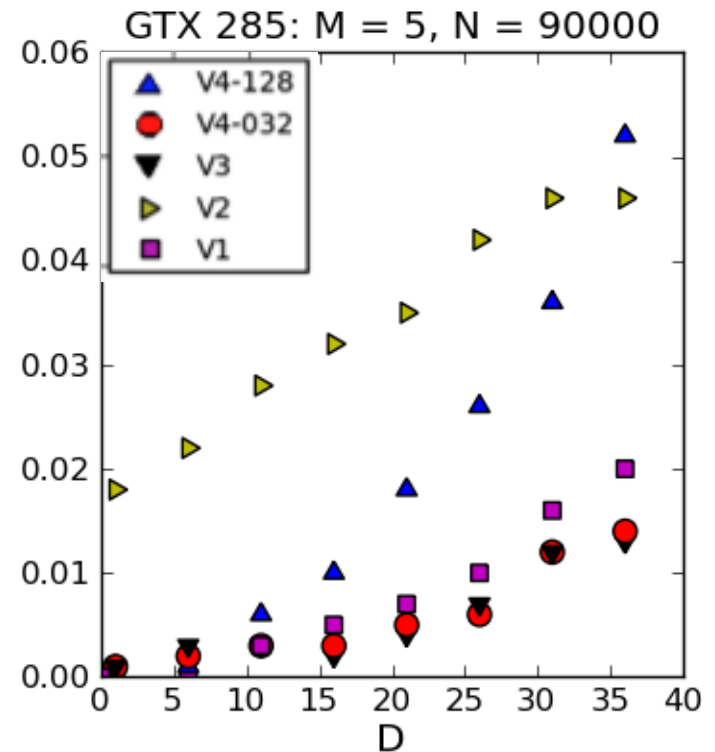
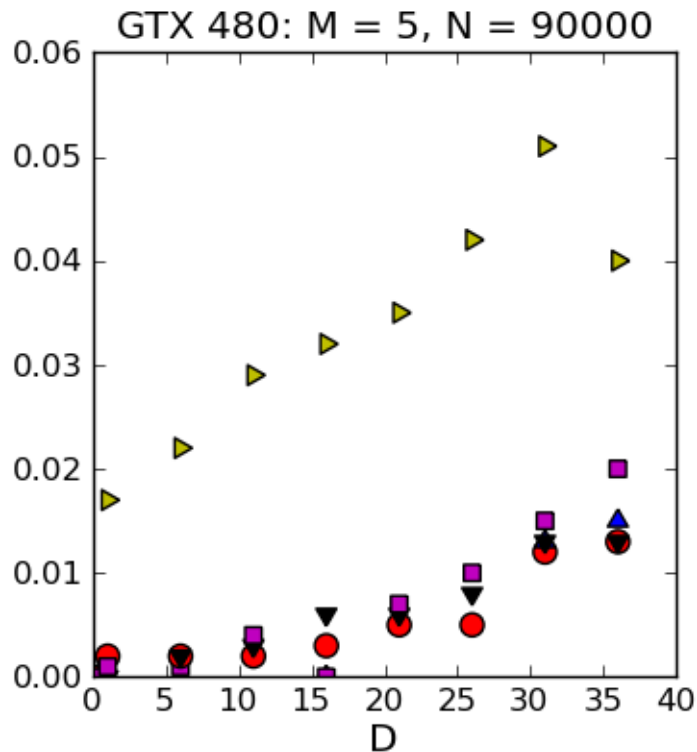
GTX₄₈₀ – Varying D

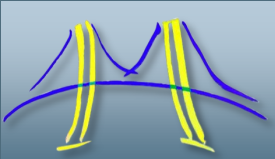




Results – Version Comparison (Raw CUDA)

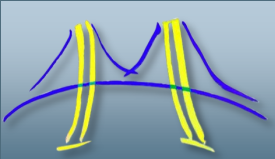
GTX285 vs. 480



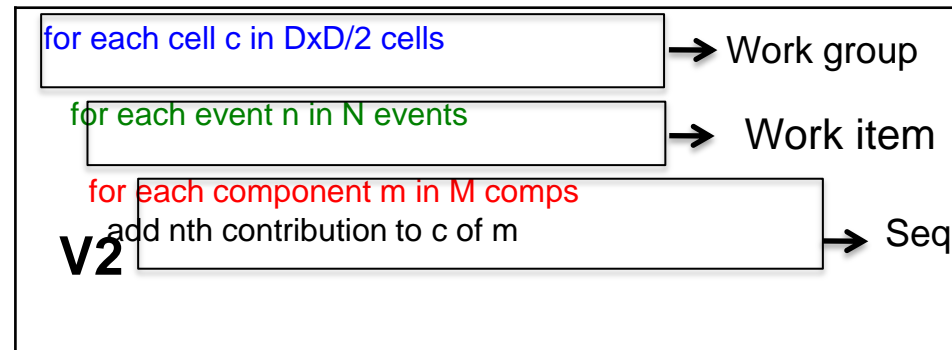
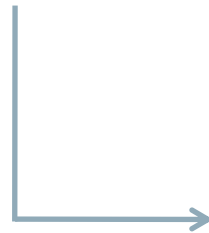
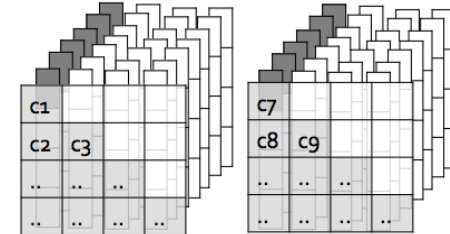
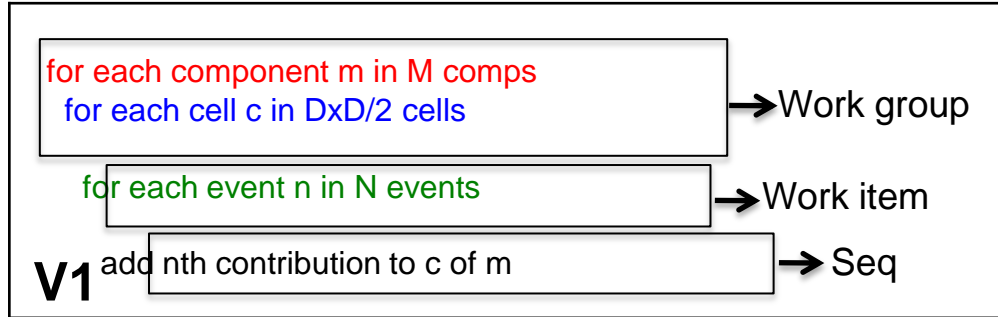


SEJITS Framework: Current Implementation

- ASP framework
 - C and CUDA compiling with CodePy (using PyCuda)
 - PyUBLAS to eliminate copies between C and Python
 - Version selection based on previous timings
- Evaluation platforms:
 - GTX480 (Fermi)
 - 14 SM, 32 SIMD, 48K shared mem, 3GB DRAM
 - GTX 285
 - 30 SM, 8 SIMD, 16K shared mem, 1GB DRAM
 - CUDA SDK 3.2
 - NVCC 3.2



Covariance Matrix Computation – Code Variants



.....