

Comparing and Combining Profiles

Serap Savari

SAVARI@RESEARCH.BELL-LABS.COM

Cliff Young

CYOUNG@RESEARCH.BELL-LABS.COM

Bell Laboratories, a division of Lucent Technologies

Rooms 2C-451 and 2C-523, 700 Mountain Avenue, Murray Hill, NJ 07974, U.S.A.

Abstract

How much do two profiles of the same program differ? When has a profile changed enough to warrant reexamination of the profiled program? And how should two or more profiles be combined to make a better hybrid profile? To answer these questions, we borrow concepts from information theory that measure the differences between the distributions of random variables. Treating the distribution of execution frequencies in a profile like the distribution of a random variable lets us use metrics of difference between distributions and induces a mathematically well-founded way to combine distributions. For concreteness, we report results based on branch profiles, but our algorithms apply to any frequency-based profiling method. Our metric agrees in character with prior metrics while allowing us to adjust the relative importance of coverage and conflict issues. Our hybrid distributions do better than previous profile blending methods.

1. Introduction

Many recent systems use profiles to guide optimization and translation towards a variety of goals. These systems have been striking in their ability to instrument with extremely low performance overhead (e.g., Compaq's Continuous Profiling Infrastructure, Anderson et al., 1997), their innovative application of profiling to binary translation (e.g., the FX!32 x86-to-Alpha translation system, Hookway & Herdeg, 1997), and their exhibition of speedups in dynamic optimization (e.g., Hewlett-Packard's Dynamo dynamic optimization project, Bala, Duesterwald, & Banerjia, 1999). All of these recent systems blur the traditional barrier between static (compiler) and dynamic (hardware) approaches. What used to be done once at compile time can now be done (and redone) at run time. But to do so effectively, we need to be able to compare profiles, to determine when the behavior of a program has changed significantly, and to combine multiple profiles to produce a better hybrid profile. This work uses information theoretic ideas to attack these problems.

Information theory deals with efficient representations of data streams and related issues like coding, compression, and prediction. Information theory is a related field to computer science and statistics, so it is not surprising that ideas from information theory are relevant to problems in profiling. Some information theory has been used recently to analyze dynamic branch prediction schemes (Chen, Coffey, & Mudge, 1996, and Federovsky, Feder, & Weiss, 1998). To our knowledge, this is the first work that applies information theoretic concepts to profiling.

Profiles are statistics about the execution of a program; they are commonly execution frequencies but they can also include data about system performance at a variety of levels. Our information-theoretic approach applies to any frequency-based profiling method. For concreteness and

simplicity of exposition, we constrain our discussion in this paper to branch profiles; much of the relevant scholarly work also relates to branch profiling.

Our first two problems, comparing profiles and determining when program behavior has changed, suggest that we need some sort of metric of the differences between profiles. Prior work in this area includes the concepts of coverage, conflict, and a geometric measure introduced recently by Kistler and Franz (1998). In their extensive cross-validation study of static branch prediction, Fisher and Freudenberger discuss *coverage*, noting that some runs of a program execute different parts of the program from other runs (Fisher & Freudenberger, 1992). They experimented with a number of different metrics but did not find any one that satisfactorily explained their results. Calder, Grunwald, and Srivastava later defined coverage as the percentage of branches from a program run that were also executed during a profile training run (Calder, Grunwald, & Srivastava, 1995). We adopt the latter definition of coverage for this paper, and we define both static coverage and dynamic coverage. *Static coverage* is the number of branch instructions touched in both the training and testing runs divided by the number of branch instructions touched during the testing run. *Dynamic coverage* is similar, but it weighs each branch by its execution frequency. The dynamic coverage is the sum of execution frequencies of instructions in the testing run that were also touched during the training run divided by the total branch execution frequency of the testing run.

Conflict is a relative of coverage: it measures the percentage of branches that change majority direction between the training and testing data sets. Similarly to coverage, we also define statically and dynamically weighted forms of conflict. Coverage and conflict are related in the way that control-flow graph (CFG) node profiles and CFG edge profiles are related: both give frequencies for points of the program, but conflict metrics and edge profiles also capture the exit bias of individual branch instructions. For reference purposes, Table 1 summarizes the definitions of the various kinds of coverage and conflict. We will return to Kistler and Franz’ metric later in the paper; it is better to discuss it after we have introduced our own metric.

Static Coverage	% of <i>test</i> instructions also touched during <i>train</i>
Dynamic Coverage	% of executed instructions on <i>test</i> that were also executed during <i>train</i>
Static Conflict	% of <i>test</i> branch instructions whose bias disagrees with the <i>train</i> bias
Dynamic Conflict	% of executed branch instructions on <i>test</i> whose bias disagrees with the <i>train</i> bias

Table 1: Basic profile metrics.

Fisher and Freudenberger also examined our third problem, combining multiple profiles into a hybrid profile. They investigated different ways of combining profiles that they called *unscaled*, *scaled*, and *polling*. Unscaled combination simply adds the profile counts from different runs; longer runs therefore carry more weight. Scaled combination divides each data point frequency by the total frequency for that profile, then arithmetically averages the weights seen in different profile runs. Polling gives each data set one vote no matter how large the program run. Fisher and Freudenberger report that polling performs poorly, while scaled and unscaled combining both give good results; they report scaled results because scaling seems intuitively better.

So far, we have been rather vague about what we mean by a, “better,” method to build hybrid profiles. Each designer has a set of goals, including but not limited to power, performance, cost, and so forth; a better profile would help with respect to such goals. In addition, certain domain-

specific knowledge may change which profiles are better compared to others. For example, one might know that recent profiles are more accurate than older profiles for a particular application, or one might know that longer profile runs correspond to more common uses of a program. Each of these specific pieces of knowledge should be used to build a better-performing system *if it is available*. Our approach addresses the case where one has collected a set of profiles but has no additional external knowledge about what should or should not be valued. And we also believe that our approach would still work for combining profiles in the presence of additional knowledge; we would just need to bias some of the parameters to favor the more likely cases.

The next section introduces our metric, which is based on the information-theoretic concept of relative entropy or Kullback-Liebler distance. Section 3 then gives an algorithm using geometric averaging and relative entropy to blend a pair of profiles. We present empirical results in Section 4. Lastly, we discuss the implications and usefulness of this work.

2. Relative Entropy

Our approach treats a frequency profile as a probability distribution. This makes intuitive sense: a profile says where a program is likely to spend its time, while a probability distribution says which values a random variable is likely to assume. In effect, we treat the program counter as a random variable (even though we know it is not—more about this later) then view the frequency profile as a set of samples drawn from the program counter’s underlying probability distribution. The statistical concept of a probability distribution is like the profiling ideal of program behavior. We can never exactly know either, but we can approximate them with increasing accuracy as we see more samples or more behavior.

We translate from profile frequencies to probabilities in the obvious way: scale the frequencies by the total frequency over the profile. This is just the first step of the scaled combination method. More formally, let S denote the set of branch instructions in a program and consider the branch trace of the program, which is a sequence of ordered pairs $(s, j) \in S \times \{0, 1\}$. An appearance of $(s, 0)$ indicates that an execution of the program entered branch s and fell through, while $(s, 1)$ indicates that the branch s was taken. Let $N_{s,j}, (s, j) \in S \times \{0, 1\}$ denote the number of times ordered pair (s, j) was encountered in the trace. We call our translation the *independent and identically distributed (i.i.d.) branch and decision model*. In this model, we assume that the sequence of ordered pairs is i.i.d.¹ with probability

$$p_{s,j} = \frac{N_{s,j}}{\sum_{t \in S, i \in \{0, 1\}} N_{t,i}}, (s, j) \in S \times \{0, 1\}.$$

1. Independent and identically distributed is a term from introductory probability. “Identically distributed,” means that each value is drawn from the same underlying probability distribution. “Independent,” means that the value observed for one sample does not affect the value of other samples.

At the risk of illustrating the obvious, Table 2 lists some sample frequency profiles and their translations into scaled probabilities under the i.i.d. branch and decision model. Each row of the table shows values for one profile. Each column shows the frequencies or probabilities associated with a single event in the profiled program; an event in an edge profile would be an ordered pair $(s, j) \in S \times \{0, 1\}$. Note that *prof1* covers all of the events, that *prof2* and *prof3* are very similar, that *prof5* covers the same events as *prof2* and *prof3* but with different emphasis, while *prof4* overlaps very little with *prof2*, *prof3*, and *prof5*. We will use these profiles to help explain this section and the next section.

	Profile	<i>event1</i>	<i>event2</i>	<i>event3</i>	<i>event4</i>
Frequencies	<i>prof1</i>	250	250	250	250
	<i>prof2</i>	1500	500	20	0
	<i>prof3</i>	300	150	10	0
	<i>prof4</i>	0	0	400	450
	<i>prof5</i>	400	500	600	0
Probabilities	<i>prof1</i>	0.25	0.25	0.25	0.25
	<i>prof2</i>	0.74	0.25	0.01	0.00
	<i>prof3</i>	0.65	0.33	0.02	0.00
	<i>prof4</i>	0.00	0.00	0.47	0.53
	<i>prof5</i>	0.27	0.33	0.40	0.00

Table 2: Sample profile frequencies and their translation into probabilities under the i.i.d branch and decision model.

Before we can discuss the relative entropy of two probability mass functions, we need to define the *entropy* of a single random variable. Let X be a discrete random variable taking values from an alphabet Ψ with probability mass function $p(x) = \Pr\{X = x\}$, $x \in \Psi$. Entropy measures the uncertainty of a random variable; it is a lower bound on the average length of the shortest description of X (see Cover & Thomas, 1991, section 1.1). Said differently, entropy is the average number of bits required per symbol to optimally encode a stream of symbols chosen from Ψ using the probability distribution p . The entropy $H(X)$ of X is the expected value of $\log(1/p(x))$, where X is drawn according to $p(x)$, i.e.,

$$H(X) = \sum_{x \in \Psi} p(x) \log \frac{1}{p(x)}$$

Because of continuity arguments, we use the convention that $-0 \log 0 = 0$. In encoding a stream of values, we use more bits to represent the rarer symbols; the log of the inverse of the probability of a symbol is the right number of bits to use to encode that symbol. To relate this definition of entropy back to dynamic optimization and translation, the entropy of a profile distribution can be seen as a measure of locality in program execution; a program that spends most of its time in a few places (large $p(x)$ at few values of x) will have low entropy, while a program with poor locality

(small $p(x)$ at many values of x) will have high entropy. Dynamic optimization and translation systems exploit the principle of locality, so we expect such systems to do better on programs whose profiles have low entropy.

Table 3 shows the entropy of each sample profile and the components that each event contributes to the total. Each component is a term of the form $p(x)\log(1/p(x))$, where x is just one of the profiled events. Low-probability events contribute significant components, so large numbers of low-probability events will raise the entropy of a profile. Profiles with high entropy tend to have large numbers of events that are evenly balanced in probability. Profiles with low entropy concentrate their frequencies and probabilities in a small number of events.

Profile	Entropy components				Entropy
	<i>event1</i>	<i>event2</i>	<i>event3</i>	<i>event4</i>	
<i>prof1</i>	0.50	0.50	0.50	0.50	2.00
<i>prof2</i>	0.32	0.50	0.07	0.00	0.88
<i>prof3</i>	0.40	0.53	0.12	0.00	1.05
<i>prof4</i>	0.00	0.00	0.51	0.49	1.00
<i>prof5</i>	0.51	0.53	0.53	0.00	1.57

Table 3: Per-event components of entropy and total entropy for each sample profile.

Entropy tells us something about the locality of a single profile run, but our goal remains to compare two different profiles. The *relative entropy* or *Kullback-Liebler distance* of two probability distributions on Ψ is a measure of their distance. Relative entropy measures the inefficiency of assuming that the distribution is q when it is actually p . This is like the difference between optimizing or translating for some assumed profile when the program run corresponds to a different actual profile. If we are given two probability mass functions $p(x)$ and $q(x)$, then their relative entropy $D(p \parallel q)$ is defined by

$$D(p \parallel q) = E_p\left(\log\frac{p(X)}{q(X)}\right) = \sum_{x \in \Psi} p(x)\log\frac{p(x)}{q(x)}.$$

Similarly to before, we use the convention that $0\log(0/q) = 0$ and $p\log(p/0) = \infty$; in a moment we will also adjust some of the profile values to eliminate the infinite terms. The description length interpretation of $D(p \parallel q)$ is that $H(p) + D(p \parallel q)$ is a lower bound on the average length of the shortest description of X when p is the true distribution of the random variable and q is the distribution for which the description was designed (see Cover & Thomas, 1991, section 2.3). In other words, $D(p \parallel q)$ is the average additional number of bits that are transmitted because of the incorrectly assumed distribution. Relative entropy measures a mismatch between probability mass functions that is like the mismatch between training and testing data sets in profile-based optimization: we optimize with a training profile (assumed distribution), but realistic results are measured with a different, testing profile (actual distribution).

Relative entropy is always non-negative (see Cover & Thomas, 1991, section 2.6) and is zero only when the two distributions are identical². Although it does provide a sense of the closeness of two distributions, relative entropy is not truly a distance metric because it is not symmetric and it does not satisfy the triangle inequality. However, the asymmetry models a realistic aspect of profiles: one profile can be a good training set for the other without the reverse being true.

Observe that the relative entropy

$$D_{edge}(p \parallel q) = \sum_{(s,j) \in S \times \{0,1\}} p_{s,j} \log \frac{p_{s,j}}{q_{s,j}}$$

is infinite when there is a term of the form $p \log(p/0)$. Similarly to the continuity arguments above, we modify any probability distribution q so that $q \geq \epsilon$ for all x while we simultaneously maintain the ratios among the remaining probability mass values in q . ϵ should be chosen smaller than the inverse of the total execution frequency; otherwise the modified values can end up larger than the scaled non-zero frequency events. The choice of ϵ indicates how much we wish to penalize a lack of coverage (smaller values of ϵ will give larger metric values for points where the training set does not cover the testing set). For the experimental results in this paper, we used a value of $\epsilon = 1 \times 10^{-30}$.

D_{edge} is a more sophisticated metric than coverage and conflict, because D_{edge} examines the relative frequencies of nodes and edges in the assumed and actual profiles rather than summarizing nodes by presence or absence and edge pairs by majority direction. We would expect that such continuous metrics would be more sensitive than binary ones: e.g., the difference between branch biases of 5% and 95% on different profiles is much larger than the differences between biases of 49% and 51%.

Table 4 shows the relative entropies of pairs of our sample profiles. The data for each actual data set appear in rows, while the assumed data sets appear in columns.³ The diagonal of the matrix is always zero: a profile exactly matches itself. We see very high relative entropies when the assumed profile does not cover the actual profile, e.g., no other profile is a good assumed data set for *prof1*, while *prof4* is a bad assumed data set for all other actual data sets. On the other hand, *prof2* and *prof3* have low relative entropy because they are close in character. *Prof5* covers the same events as *prof2* and *prof3*, but with different emphasis; this distinction is correctly reflected by the slightly higher values for relative entropy when using *prof5* rather than *prof2* or *prof3* for either the assumed or the actual distributions.

One aspect of D_{edge} is unsatisfying: nothing relates the edges that exit a single control-flow graph node to one another. I.e., the counts $N_{s,0}$ and $N_{s,1}$ are related by the logic of the program in a way that $N_{s,0}$ and $N_{t,0}$ are not, but our metric does not capture this difference. We also

2. One aspect of the formula seems unintuitive: the logarithm allows us to have negative terms when $p(x) < q(x)$. While this is true for any particular value of x , we never achieve a negative total because negative terms are offset by positive terms where the actual distribution is larger than the assumed distribution.

3. These are usually called testing data set (actual) and training data set (assumed). SPEC already uses these terms for their three input data sets: “ref”, “train,” and “test.” So we use *assumed* and *actual* to avoid confusion.

Assumed	<i>prof1</i>	<i>prof2</i>	<i>prof3</i>	<i>prof4</i>	<i>prof5</i>
<i>prof1</i>	0.00	25.19	24.85	48.33	24.12
<i>prof2</i>	1.12	0.00	0.03	97.80	0.94
<i>prof3</i>	0.95	0.03	0.00	96.47	0.74
<i>prof4</i>	1.00	54.90	54.36	0.00	52.38
<i>prof5</i>	0.43	1.88	1.35	58.66	0.00

Table 4: Relative entropies of the sample profiles.

examined a couple of metrics that strove to capture this difference. Neither alternate metric seemed very different in character from D_{edge} , so we omit them from this paper.

Code for converting frequency counts into probability distributions and for computing the Kullback-Liebler distance between two distributions appears in Figure 1.

```

double *
Freq2Prob(int nevents, ulong *freqs)
{
    ulong total = 0;
    double *results = malloc(nevents * sizeof(double));

    for(i = 0; i < nevents; i++)
        total += freqs[i];

    for(i = 0; i < nevents; i++)
        results[i] = (double)freqs[i] / total;

    return results;
}

double
KLDist(int nevents, double *actual, double *assumed)
{
    double distance = 0.0;

    for(i = 0; i < nevents; i++)
        distance += actual[i] * log2(actual[i] / assumed[i]);

    return distance;
}

```

Figure 1: Code for converting frequencies to probability distributions and for computing the Kullback-Liebler distance between two probability distributions. Note that we have omitted checks for division by zero for the sake of brevity; actual code substitutes epsilon (ϵ) as described in the paper.

3. Combining Profiles

Our algorithm for combining two frequency profiles has a simple implementation and description, but there is a more complex argument for why the algorithm produces useful results. We present the algorithm first; those who are uninterested in the mathematical justification can skim the remainder of this section.

3.1. Algorithm

To blend two profiles, we need to find a synthetic profile that is a compromise between the two. Assume our two profiles have been translated into probability distributions using the i.i.d. branch and decision model; call them P_0 and P_1 . Then we will search values of the parametric probability mass function P_λ , defined as

$$P_\lambda(x) = \frac{P_0^{1-\lambda}(x)P_1^\lambda(x)}{\sum_{a \in \Psi} P_0^{1-\lambda}(a)P_1^\lambda(a)}$$

P_λ is a geometric weighting of P_0 and P_1 ; $\lambda = 0$ gives P_0 and $\lambda = 1$ gives P_1 , while $\lambda = 1/2$ gives the geometric average of the two probability distributions.

Figure 2 illustrates how P_λ relates to the other probability distributions. The triangle represents the probability simplex, the space of probability distributions over Ψ (the simplex is typically drawn as a triangle because a space of n possible values for the random variable has only $n - 1$ degrees of freedom—the last probability must ensure that the sum of probabilities is 1). P_0

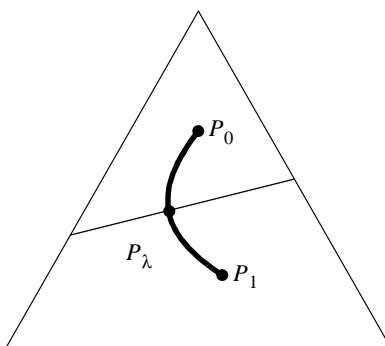


Figure 2: The probability simplex, showing the original distributions P_0 and P_1 , and the range of values of the parameterized distribution P_λ . Distributions in the upper part of the simplex are closer to P_0 , while those in the lower part of the simplex are closer to P_1 . (After Cover & Thomas, 1991, Figure 12.9).

and P_1 are points in this space. The space is divided into two parts: points that are closer to P_0 , and parts that are closer to P_1 . The line between the two spaces is analogous to the perpendicular bisector in Euclidean geometry. The values of P_λ trace out an arc from P_0 to P_1 ; we will search along this arc for our hybrid probability distribution.

Next, we use the relative entropy metric from Section 2 to choose the value of λ ; we look for a value of λ such that P_λ is equidistant from both P_0 and P_1 . We could do this directly, but there is an easier way. To do this, we build a function $c(\lambda)$ that computes

$$c(\lambda) = \sum_{x \in \Psi} P_0^{1-\lambda(x)} P_1^{\lambda(x)}$$

then binary search the interval $[0,1]$ for λ^* , the value that minimizes this function. Then P_{λ^*} is our synthesized profile distribution.

Code for computing $c(\lambda)$ and then the synthesized probability distribution λ^* appears in Figure 3.

Table 5 tabulates part of the search space when our algorithm combines *prof2* with *prof5*. The first column depicts values of λ in the interval $[0, 1]$; the second through fifth columns show the corresponding probabilities for each event and that value of λ , and the sixth column shows the value of $c(\lambda)$, which is minimized between $\lambda = 0.50$ and $\lambda = 0.75$. Note that the values for *event2* unintuitively exceed the original profile values in the range $\lambda \in [0.38, 0.88]$. This is due to the combination of geometric blending and normalization that we perform. The numerator of our formula for P_λ , the geometrically blended probabilities from the original profiles, produces terms that sum to less than 1.0 unless $P_0 = P_1$, or $\lambda = 0$, or $\lambda = 1$. To normalize the probabilities, we divide by the sum of all synthesized terms. This boosts the synthesized probabilities of events that are relatively common to both profiles for some intermediate values of λ . This boosting would never happen with arithmetic blending. The extra column in Table 5 lists the normalizing coefficients that were used. Figure 4 graphs the columns of Table 5, depicting the relative event probabilities in our synthesized profile distributions as we vary λ from 0 to 1.

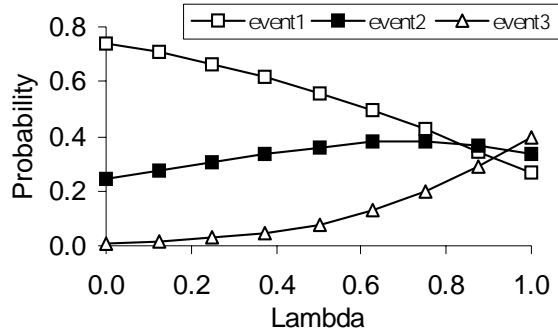


Figure 4: Synthesized probabilities as λ varies from 0 to 1.

```

double
evaluate_c(double lambda, int nevents, double *p0, double *p1)
{
    double total = 0.0;

    for(i = 0; i < nevents; i++)
        total += pow(p0[i], 1.0 - lambda) * pow(p1[i], lambda);

    return total;
}

double
search_lambda(int nevents, double *p0, double *p1)
{
    double lambdalo, lambdamid, lambdahi, distlo, distmid, disthi;

    distlo = evaluate_c((lambdalo = 0.0), nevents, p0, p1);
    disthi = evaluate_c((lambdahi = 1.0), nevents, p0, p1);
    lambdamid = (lambdalo + lambdahi) / 2.0
    distmid = evaluate_c(lambdamid, nevents, p0, p1);

    while(fabs(distlo - disthi) > STOP_BIN_SEARCH){
        if(distlo < disthi){
            lambdahi = lambdamid;
            disthi = distmid;
        }else{
            lambdalo = lambdamid;
            distlo = distmid;
        }
        lambdamid = (lambdalo + lambdahi) / 2.0
        distmid = evaluate_c(lambdamid, nevents, p0, p1);
    }

    return lambdamid;
}

double *
kl_blend(int nevents, double *p0, double *p1)
{
    double lambdastar = search_lambda(nevents, p0, p1);
    double norm = evaluate_c(lambdastar, nevents, p0, p1);
    double *results = malloc(nevents * sizeof(double));

    for(i = 0; i < nevents; i++)
        results[i] = pow(p0[i], 1.0 - lambda) * pow(p1[i], lambda) / norm;

    return results;
}

```

Figure 3: Code for evaluating the function $c(\lambda)$, for binary searching for the minimum of $c(\lambda)$ in the interval $[0,1]$, and for constructing the geometrically parameterized synthetic profile distribution.

λ	Values of P_λ				$c(\lambda)$	Normalizing coefficient
	<i>event1</i>	<i>event2</i>	<i>event3</i>	<i>event4</i>		
0.00 = <i>prof2</i>	0.74	0.25	0.01	0.00	1.00	1.00
0.13	0.71	0.28	0.02	0.00	0.93	0.93
0.25	0.66	0.31	0.03	0.00	0.87	0.87
0.38	0.62	0.34	0.05	0.00	0.82	0.82
0.50	0.56	0.36	0.08	0.00	0.80	0.80
0.63	0.50	0.38	0.13	0.00	0.79	0.79
0.75	0.42	0.38	0.20	0.00	0.81	0.81
0.88	0.35	0.37	0.29	0.00	0.88	0.88
1.00 = <i>prof5</i>	0.27	0.33	0.40	0.00	1.00	1.00

Table 5: Part of the search space of synthetic profiles when combining *prof2* and *prof5*.

Figure 5 charts the values of the relative entropies from the synthesized profile to the original profiles; it captures the trade-off from favoring one profile to the other as we vary λ from 0 to 1. λ^* is the point on the horizontal axis where the two distance curves cross: it exactly balances the distance from the synthesized distribution P_λ to each of the two original distributions.

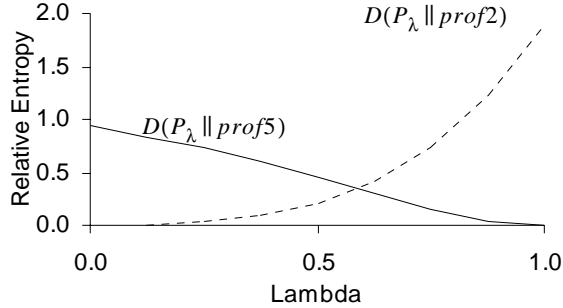


Figure 5: Values of the relative entropies from the hybrid profile P_λ to the original profiles *prof2* and *prof5*. The horizontal axis shows the values of λ . (After Cover & Thomas, 1991, Figure 12.10).

Our algorithm runs in time proportional to the number of profile events times the number of bits of accuracy desired in λ ; in all of our experiments our computations were I/O-bound.

We describe how our algorithm performs for combining static branch profiles in Section 4.2; now we continue with the formal explanation of how and why it works.

The justification for our algorithm reverses an argument from the theory of hypothesis testing. In hypothesis testing, we have a number (or in our case just a pair) of hypotheses, one of which is true. Based on some amount of evidence, we try to tell which of the hypotheses is the correct one; hypothesis testing analyzes the probability of error and allows us to weigh the hypotheses based on the evidence to minimize errors. The reversal comes from how we think about the evidence we have seen so far. Instead of using the evidence we have seen so far to guess the correct hypothesis,

we will construct a set of evidence that is simultaneously close to both original profiles but hard to use to choose which one is correct. We will start by introducing the usual terminology that describes the hypothesis testing problem, then we will consider the sorts of evidence that make decisions hard.

3.2. Terminology for Hypothesis Testing

Suppose we are given independent and identically distributed random variables X_1, X_2, \dots, X_n that have probability mass function $Q(x)$. When we observe the values of these random variables we see the sequence of specific values x_1, x_2, \dots, x_n . These values are the pieces of evidence about $Q(x)$ that we might observe. We have two hypotheses H_0 and H_1 for the nature of $Q(x)$. Either $H_0: Q = P_0$ or $H_1: Q = P_1$, i.e., H_0 states that Q , the probability distribution, is P_0 , while H_1 states that the probability distribution is P_1 . Our choice for the solution is based on a decision function $g_n(x_1, \dots, x_n)$ which returns the value zero if it chooses hypothesis H_0 and otherwise returns the value one.

With the decision rule g_n there are two error probabilities:

$$\alpha_n = \Pr(g_n(X_1, \dots, X_n) = 1 | H_0 \text{ true})$$

$$\beta_n = \Pr(g_n(X_1, \dots, X_n) = 0 | H_1 \text{ true})$$

In other words, α_n is the error probability that after seeing the first n samples, g_n incorrectly picks H_1 , while β_n is the error probability of the reverse failure. We would like to minimize both error probabilities, but there is a trade-off between them.⁴

It is possible to show that the best decision rules are likelihood tests of the form

$$\text{Choose } H_0 \text{ if } \frac{P_0(x_1, \dots, x_n)}{P_1(x_1, \dots, x_n)} > T \text{ and } H_1 \text{ otherwise}$$

for some $T \geq 0$ (see Cover & Thomas, 1991, section 12.7). In other words, consider the ratio of probabilities of seeing the empirical evidence so far under the two hypothesized probability distributions. If this ratio exceeds the threshold value T , choose H_0 ; otherwise choose H_1 . Choosing $T = 1$ turns out to give us equal probabilities of error α_n and β_n ; other values of T will favor one kind of error over the other. This kind of decision rule makes intuitive sense: it tracks which of the hypothesized probability distributions is closer to the empirical samples we have seen so far.

Let P_{X^n} be the empirical probability distribution corresponding to (x_1, \dots, x_n) . In other words, P_{X^n} is just the probability distribution constructed from the frequency of occurrences of

4. We can trivially minimize α_n or β_n by using the constant decision functions $g_n = 0$ or $g_n = 1$, respectively, but finding a decision function that makes a good trade-off between α_n and β_n requires more analysis.

different values in the series (x_1, \dots, x_n) . Then we can further show that the likelihood test

$$\frac{P_0(x_1, \dots, x_n)}{P_1(x_1, \dots, x_n)} > T$$

is equivalent to

$$D(P_{X^n} \parallel P_1) - D(P_{X^n} \parallel P_0) > \frac{1}{n} \log T.$$

(see Cover & Thomas, 1991, Section 12.7, p. 307; the argument basically takes the logarithm of the likelihood test inequality).

This concludes our definition of terms for the hypothesis testing problem. The next subsection continues with how to find the empirical distribution that corresponds to our synthesized profile.

3.3. Boundary Evidence

Using the terminology from the last subsection, we can now state our goal more precisely: find an empirical probability distribution, P_{X^n} , that is close to both P_0 and P_1 while still making the decision between them hard. Our metric tells us the distance from P_{X^n} , the distribution we have seen so far, to the two profiles, P_0 and P_1 . There is a space of such empirical distributions P_{X^n} (depicted in Figure 2), and we can use the distance metric to partition them into three pieces: distributions that are closer to P_0 , those that are closer to P_1 , and those that are equidistant to both P_0 and P_1 . Our desired distribution will be in the equidistant set, and it must also be close to the originals.

In the Bayesian approach to hypothesis testing, we assign prior probabilities to the two hypotheses, then we analyze the total error possible in the system. Call these *a priori* probabilities $\pi_j, j \in \{0, 1\}$, for each hypothesis H_j . In other words, before we see any of the evidence from the X_i , there is a π_j chance that hypothesis H_j is correct (of course, $\pi_0 + \pi_1 = 1$). Using our definitions, the *a posteriori* probability of error (the odds that we misidentify after having seen n samples) is

$$P_e^{(n)} = \pi_0 \alpha_n + \pi_1 \beta_n$$

Next, define

$$D^* = \lim_{n \rightarrow \infty} \min_{g_n} -\frac{1}{n} \log P_e^{(n)}$$

As we increase the number of samples that we see, our likelihood of making an error using the best decision rule decreases exponentially. Roughly speaking, D^* is the best rate of decay or base of this exponential curve; it describes the rate at which errors decrease as we see more samples using the best possible decision rule. D^* exists and is well-defined. The proof is a significant information theoretic result; the result is known as the Chernoff Bound. Cover derives the Chernoff Bound in Section 12.9 (Cover & Thomas, 1991); explaining the details of the bound is unfortunately beyond the scope of this work.

We have carefully chosen our definitions to match the conditions of the Chernoff Bound; it tells us that D^* turns out to be the best achievable error exponent; and further that

$$D^* = D(P_{\lambda^*} \| P_0) = D(P_{\lambda^*} \| P_1)$$

and

$$D^* = -\log \sum_{x \in \Psi} P_0^{1-\lambda^*}(x) P_1^{\lambda^*}(x)$$

where

$$\lambda^* = \arg \min_{0 \leq \lambda \leq 1} \sum_{x \in \Psi} P_0^{1-\lambda}(x) P_1^{\lambda}(x)$$

i.e., λ^* is the value of λ that minimizes the maximum of $D(P_\lambda \| P_0)$ and $D(P_\lambda \| P_1)$, the relative entropy distances from P_λ to P_0 and P_1 . Further, the function that λ^* minimizes (which is just $c(\lambda)$ from Section 3.1) is convex in the range $[0,1]$, so λ^* exists and is unique.

More slowly and working backwards, the third equation above tells us that λ^* minimizes $c(\lambda)$, the same function we minimized in the algorithm of Figure 3. Then plugging the definition of $c(\lambda^*)$ into the second equation gives us $D^* = -\log c(\lambda^*)$. Lastly, the first equation tells us that the point (λ^*, D^*) is the point where the curves cross in Figure 5.

The value of λ^* above is exactly the λ^* that we found using the algorithm of Section 3.1. So the Chernoff bound tells us that we have indeed found the boundary piece of evidence that is simultaneously close to both hypothetical distributions. When we use this to build a synthesized profile distribution, we get the profile distribution that is equidistant from both training (assumed) data sets while simultaneously also being as close as possible to both data sets.

4. Experimental Results

We collected profiles on a 533MHz Alpha 21164 clone with a 164LX motherboard, 2MB of board cache, and 128MB of main memory. The Alpha runs Digital UNIX version 4.0D (revision 878); we profiled using Atom version 2.47 and custom analysis files. We compiled the SPECint95 benchmarks using base optimization settings. Due to current limitations of our analysis code, we collected statistics for only the last run of each training, testing, or reference input set. Table 6 summarizes the number of static and dynamic branches our profiles recorded on each profiling run. We use static and dynamic in the same sense as in the definitions for our metrics: one static branch instruction in the program text might be executed many times during a run; each of these executions of the static branch is a dynamic branch.

Training run	All	ref		test		train	
Benchmark	Static	Static	Dynamic	Static	Dynamic	Static	Dynamic
099.go	5740	5603	3669393694	5464	1833855003	5046	59750628
124.m88ksim	1418	1408	8135706733	1219	61291271	1164	11682553
126.gcc	19370	14592	26190611	18565	199369904	18644	202317776
129.compress	430	430	4975851051	407	216427	415	3496670
130.li	879	870	8605396422	635	148521036	636	25032723
132.ijpeg	1373	1338	1515921676	1354	42129447	1365	93933314
134.perl	2307	2080	3114254002	1630	1186183	2047	5141087
147.vortex	6501	6482	9334649279	6488	1088823792	6487	302461902

Table 6: Benchmark and profile branch counts.

We collected three profiles for each benchmark, one for each of the reference, testing, and training inputs. These are described by the file SPEC95/doc/RUN.txt, which is part of the SPEC95 distribution:

...“ref” is the input set for generating publishable numbers, “train” is used for generating feedback directed optimization (if used) and “test” is a short input for helping verify that the benchmark was compiled correctly.

When reporting SPEC results, only the “train” inputs can be used for profiling. Typically the train and test inputs are smaller than the reference inputs.

4.1. Metrics for Comparing Profiles

Evaluating the effectiveness of a metric is difficult: a “good” metric is one that turns out to predict or measure something useful. We do not have a dynamic profiling/optimization system in which to evaluate our metric, so this section uses static branch prediction accuracy as a stand-in for performance. This is somewhat unsatisfying: we would really like to try our metric in an operating dynamic profiling and optimization system with a concrete goal such as performance, power, cost, etc. Evaluating our metric in other systems with other performance criteria is beyond the scope of this study.

Table 7 lists a plethora of statistics for each benchmark. The columns list assumed-actual profile pairs, i.e., each statistic was collected by training with the assumed data set then testing using the actual data set. There are seven rows for each benchmark. The first four show static and dynamic values for the coverage and conflict metrics defined in Section 1; these have values between 0 and 1. To simplify viewing, we have shaded all cases where dynamic coverage was less than 95% or dynamic conflict exceeded 5%. The next row lists our edge-based relative entropy metric; this ranges from 0 to infinity. The last two rows show the static branch prediction accuracies; these take values from 0 (bad) to 1 (perfect). The “resub.” row lists the prediction accuracy from training and testing on the actual data set (it ignores the assumed data set, so we see identical numbers for the same actual data set); this is provided as an upper bound on possible static prediction accuracy. The “x-validated” row shows the prediction accuracy when training on the assumed data set and running on the actual data set; in all cases the cross-validated value will be less than the resubstitution value. To simplify viewing, we have shaded all cases where prediction accuracy drops by more than 2% from the resubstitution value to the cross-validated value.

Many of the static coverage values seem low: every benchmark except *132.jpeg* and *147.vortex* have some assumed-actual pair where the static coverage is lower than 95%. In contrast, only a few of the dynamic coverage statistics are poor (*124.m88ksim* test-*, *129.compress* test-ref, four of the 6 cases on *134.perl*). For these profiles, low dynamic coverage seems a much better indicator of poor prediction performance than low static coverage. Unsurprisingly, low dynamic coverage always occurs with a large drop from resubstitution prediction accuracy to cross-validated prediction accuracy (when there is no coverage from the assumed data set, we predict branches will be taken). Coverage is a necessary but not sufficient condition for a good profile.

Static conflict is as unilluminating as static coverage: all benchmarks except *129.compress*, *132.jpeg*, and *147.vortex* have some assumed-actual pair with greater than 5% static conflict. Dynamic conflict is much better: all values greater than 5% correspond to cases where prediction accuracy has dropped by more than 2% from resubstitution to cross-validation. This is also unsurprising, as dynamic conflict is almost the definition of cross-validated misprediction rate.⁵ Using high dynamic conflict as an indicator that profiles differ misses a few cases: 3 cases on *124.m88ksim*, 3 cases on *129.compress*, and 4 cases on *130.li*. Tightening the threshold to 2% would catch all of these cases without generating any false positives.

Our metrics based on relative entropy do a reasonable but imperfect job identifying cases where prediction accuracy drops from resubstitution to cross-validation. For most of the bad prediction accuracy cases, the metrics attain values above 1. But this threshold gives some false negatives: *129.compress* ref-train, two cases on *130.li*, and there are also false positives: *124.m88ksim* ref-train, *126.gcc* ref-test, *129.compress* train-test, and *132.jpeg* ref-train and ref-test. Using a different threshold than 1 trades off false positives for false negatives. These results are encouraging because they suggest that our metrics are indeed useful for measuring profile similarity; they would be more encouraging if we could beat dynamic conflict as well. But as we see in the next

5. Most of the numbers appear as if the dynamic coverage is the difference between resubstitution and cross-validation prediction accuracy. This would only be the case if branches that reversed direction always fell through or always jumped. We can infer that most of the execution frequency of branches that reversed direction came from strongly biased branch instructions.

Assumed	test-	train-	ref-	train-	ref-	test-	test-	train-	ref-	train-	ref-	test-
Actual	ref	ref	test	test	train	train	ref	ref	test	test	train	train
Metric												
stat. cover	0.9522	0.8938	0.9764	0.8840	0.9925	0.9572	0.7195	0.7310	0.9858	0.9150	1.0000	0.9135
dyn. cover	1.0000	0.9998	0.9986	0.9984	1.0000	1.0000	0.9883	0.9827	0.9941	0.9793	1.0000	0.9982
stat. conf.	0.1087	0.1624	0.0943	0.1720	0.1140	0.1375	0.2379	0.2276	0.0882	0.1008	0.0645	0.0912
dyn. conf.	0.0025	0.0074	0.0030	0.0086	0.0050	0.0059	0.0583	0.0418	0.0780	0.0374	0.0395	0.0464
KL edge	0.0316	0.2792	0.1503	0.3771	0.1982	0.1887	1.4965	2.4849	1.0253	2.7121	0.4329	0.8071
resub.	0.7709	0.7709	0.7723	0.7723	0.7555	0.7555	0.8512	0.8512	0.8751	0.8751	0.8800	0.8800
x-validated	0.7684	0.7635	0.7693	0.7637	0.7505	0.7496	0.7929	0.8094	0.7971	0.8377	0.8405	0.8336
stat. cover	0.8636	0.8217	0.9975	0.8154	0.9940	0.8540	0.9963	0.9948	0.9845	0.9963	0.9751	0.9883
dyn. cover	0.9113	0.9969	1.0000	0.9950	1.0000	0.8498	1.0000	1.0000	0.9816	1.0000	0.9866	0.9999
stat. conf.	0.0902	0.1321	0.0353	0.1526	0.0395	0.1091	0.0120	0.0105	0.0199	0.0126	0.0249	0.0190
dyn. conf.	0.0238	0.0207	0.0003	0.0446	0.0180	0.0589	0.0001	0.0001	0.0096	0.0008	0.0049	0.0004
KL edge	8.9533	3.7468	0.2066	5.4362	2.0927	17.3343	0.4944	0.1775	2.6925	0.1655	1.5523	0.1590
resub.	0.9209	0.9209	0.9291	0.9291	0.9373	0.9373	0.8909	0.8909	0.8902	0.8902	0.8954	0.8954
x-validated	0.8971	0.9002	0.9288	0.8845	0.9192	0.8784	0.8908	0.8908	0.8806	0.8893	0.8905	0.8950
stat. cover	0.9885	0.9908	0.7769	0.9624	0.7755	0.9583	0.6750	0.9837	0.8613	0.8601	0.9995	0.6849
dyn. cover	0.9999	0.9999	0.9860	0.9997	0.9964	0.9999	0.8935	1.0000	0.7293	0.7293	1.0000	0.8814
stat. conf.	0.0539	0.0527	0.1671	0.0617	0.1658	0.0630	0.2375	0.0212	0.1564	0.1571	0.0142	0.2325
dyn. conf.	0.0023	0.0037	0.0060	0.0031	0.0052	0.0043	0.0901	0.0013	0.2263	0.2277	0.0015	0.0859
KL edge	0.1553	0.1536	1.3982	0.1939	0.4707	0.1516	13.4143	0.0875	34.2514	32.8709	0.1551	14.1100
resub.	0.8911	0.8911	0.8892	0.8892	0.8759	0.8759	0.9168	0.9168	0.9370	0.9370	0.9126	0.9126
x-validated	0.8888	0.8874	0.8832	0.8860	0.8708	0.8716	0.8266	0.9155	0.7107	0.7093	0.9112	0.8267
stat. cover	0.9465	0.9651	1.0000	1.0000	1.0000	0.9807	0.9981	0.9978	0.9972	0.9997	0.9971	0.9998
dyn. cover	0.7824	0.9990	1.0000	1.0000	1.0000	0.9913	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
stat. conf.	0.0465	0.0395	0.0246	0.0147	0.0241	0.0217	0.0068	0.0120	0.0079	0.0076	0.0129	0.0076
dyn. conf.	0.2142	0.1376	0.0539	0.0121	0.0511	0.0383	0.0002	0.0007	0.0006	0.0004	0.0006	0.0004
KL edge	24.0538	1.5464	6.1607	2.1490	0.6366	2.5058	0.2340	0.7902	0.1591	0.1149	0.3289	0.0726
resub.	0.8681	0.8681	0.9412	0.9412	0.8484	0.8484	0.9857	0.9857	0.9802	0.9802	0.9782	0.9782
x-validated	0.6539	0.7305	0.8873	0.9292	0.7972	0.8100	0.9855	0.9850	0.9796	0.9798	0.9776	0.9778

Table 7: Metrics comparing pairs of profile inputs

section, our metrics are beneficial when synthesizing a hybrid profile from a pair of original profiles.

4.2. Combining Profiles Using Relative Entropy

Table 8 shows prediction accuracies for a variety of combining strategies for each benchmark. Each column lists the actual (testing) profile; in that column, the other two profiles were combined to build the training profile. As in Table 7, the “resubstitution” row shows the upper bound prediction accuracy from training and testing on the same benchmark. The rows labeled “unscaled,” “scaled,” and “polling” correspond to Fisher and Freudenberger’s unscaled, scaled, and polling

Combiner	Assumed	test	ref	ref	test	ref	ref	
	Actual	train	train	test	train	train	test	
resubstitution	099.go	0.7709	0.7723	0.7555	130.li	0.8512	0.8751	0.8800
unscaled		0.7683	0.7693	0.7503		0.8016	0.7971	0.8403
scaled		0.7683	0.7683	0.7504		0.8015	0.8282	0.8618
polling		0.7653	0.7671	0.7503		0.7975	0.7864	0.8252
KL		0.7669	0.7659	0.7504		0.8136	0.8486	0.8732
resubstitution	124.m8ksim	0.9209	0.9291	0.9373	132.ijpeg	0.8909	0.8902	0.8954
unscaled		0.9204	0.9288	0.9192		0.8908	0.8890	0.8949
scaled		0.9206	0.9204	0.9190		0.8908	0.8894	0.8914
polling		0.8955	0.9259	0.8787		0.8909	0.8806	0.8905
KL		0.9200	0.9288	0.9192		0.8908	0.8893	0.8950
resubstitution	126.gcc	0.8911	0.8892	0.8759	134.perl	0.9168	0.9370	0.9126
unscaled		0.8886	0.8862	0.8717		0.9147	0.7107	0.9112
scaled		0.8886	0.8866	0.8720		0.9126	0.7052	0.9084
polling		0.8876	0.8831	0.8700		0.8497	0.7093	0.8485
KL		0.8889	0.8867	0.8720		0.9058	0.7107	0.9112
resubstitution	129.compress	0.8681	0.9412	0.8484	147.vortex	0.9857	0.9802	0.9782
unscaled		0.7305	0.8873	0.7972		0.9855	0.9796	0.9776
scaled		0.7300	0.9118	0.8013		0.9851	0.9800	0.9778
polling		0.6496	0.9176	0.8003		0.9851	0.9798	0.9777
KL		0.7305	0.8873	0.8253		0.9850	0.9798	0.9778

Table 8: Prediction accuracies for combining methods for all benchmarks

methods, respectively. The row labeled, “KL” gives the prediction accuracy from building a hybrid table using the parameterized distribution P_λ then minimizing the relative entropies from P_λ to the two training data sets. For each benchmark and each column, we have shaded the best (largest) non-resubstitution prediction accuracy entries.

Two profiles are usually better than one: the best blending method often does better than the best cross-validated single-profile prediction accuracy (the exceptions are *099.go* test-ref and ref-train, *129.compress* train-test, and *134.perl* train-ref). This matches our intuitions that blending multiple profiles ought to give a better description of program performance than using a single profile will.

In the contest of combining methods, the unscaled method is best 9 times out of 24; the scaled method is best 7 times; the polling method is best twice (matching Fisher and Freudenberger’s observations), and our “KL” method is best 15 times (counting ties for all methods that tie). The “KL” method is also very close when it is not the absolute best; the only bad cases are *129.compress* test and *134.perl* ref.

A slightly more quantitative metric is the difference between each combining method and the empirically best combining method for each benchmark. In the learning theory community, this difference is known as *regret*. More formally, if we use methods a , b , and c on a test and they achieve scores S_a , S_b , and S_c respectively, with S_c the best, then the regret for method a is $S_c - S_a$,

the regret for method b is $S_c - S_b$, and the regret for method c is 0. The average regret works out to 1.97% for the unscaled method, 1.23% for the scaled method, 5.70% for the polling method, and 0.57% for our “KL” method.

The differences in prediction accuracies listed in Table 8 are not often very large and the trends are noisy and hard to detect. For example, the “KL” combining method does significantly better than the scaled method only on the benchmark *130.li*; the other benchmarks and data sets show much smaller differences. Nevertheless, for these benchmarks and inputs, our relative-entropy-based combining method turns out well. It seems worthwhile to look at other performance metrics, benchmarks, and data sets; we expect our techniques to work well in other contexts.

Unfortunately, the SPEC rules do not currently allow blending the “train” and “test” input data sets to use on the “ref” input data set. However, for the cases where there are multiple program runs within an input data set, our combining method may do a better job of combining the profiles from those multiple runs.

5. Discussion

Kistler and Franz (1998) treated each profile as a vector in the space \mathbf{R}^n , where n is the number of items in the profile (nodes, edges, or paths). They then built a metric with values in the range [0,1] based on the geometric angle and the vector distance between the two vectors. Their metric is symmetric: the distance from profile A to profile B is the same as the distance from profile B to profile A. This does not reflect the way that profiles can be useful: profile A may be a good training set for profile B without the reverse being true. Kistler and Franz’ goal was to determine when program behavior had changed significantly enough to warrant reoptimization; their metric may be very good at this but they unfortunately do not include any experimental examples.

One obvious problem with our metric and our combining method is their heavy reliance on floating point computations. It is future work to find fast, on-line ways to compute them using simpler operators and/or integers. Another obvious practical problem is how to combine more than two profiles at the same time.

An intriguing work by Wang and Rubin suggests that it is possible to combine too many profiles and over generalize beyond the point that is useful to a particular user (Wang & Rubin, 1998). In other words, combining a few profiles turned out to be beneficial, but after a certain point, additional profiles led to decreasing and eventually negative marginal benefits from an individual user’s perspective. It would be interesting to see if our metrics can shed light on the nature of this problem.

Point profiles make an independence assumption in summarizing frequencies of profiled events without regard to the correlations between the events. Our relative entropy metric, D_{edge} , makes a similar simplifying assumption in stating that the events used to make up the profiles are independent and identically distributed. But events in the program trace are not independent or identically distributed: seeing one node or edge strongly constrains the next node or edge that will appear. Path profiles attempt to capture the correlations between some limited number of events in the program; relative entropy similarly generalizes into a Markov-chain based metric called divergence; the size of the Markov chain is similar to the profiling depth in a general path profile (Young, 1997). We have not yet applied divergence to path profiles, but we expect to be able to compare and combine them in a similar manner to what we have shown here.

6. Conclusion

We introduced a new approach to analyzing, comparing, and combining profiles that uses information theoretic ideas. Our comparison metric appears useful but would be more interesting if we could test it in a real system. Our combining method beats all previous approaches to combining profiles from different runs of our programs. There are many concepts in information theory that directly relate to problems in feedback-directed optimization; we expect to find better ways to analyze and explain the behavior of such systems and better and new techniques through the fusion of the two fields.

Acknowledgements

We thank the anonymous reviewers and Michael D. Smith for useful suggestions about how to make the paper more readable.

References

- Anderson, J., Berc, L., Dean, J., Ghemawat, S., Henzinger, M., Leung, S., Sites, R., Vandevoorde, M., Waldsburger, C., & Weihl, W. (1997). Continuous Profiling: Where Have All the Cycles Gone? *ACM Transactions on Computing Systems*, 15(4):357-390.
- Bala, V., Duesterwald, E., & Banerjia, S. (1999). Transparent Dynamic Optimization: The Design and Implementation of Dynamo. Tech. rep. HPL-1999-78, HP Laboratories.
- Calder, B., Grunwald, D., & Srivastava, A. (1995). The Predictability of Branches in Libraries. Tech. rep. 95/6, Digital Western Research Laboratory.
- Chen, I., Coffey, J., & Mudge, T. (1996). Analysis of Branch Prediction Via Data Compression. *In Proc. Seventh Intl. Conf. on Architectural Support for Prog. Lang. and Operating Systems*. New York, NY: ACM.
- Cover, T., & Thomas, J. (1991). *Elements of Information Theory*. New York, NY: John Wiley & Sons.
- Federovsky, E., Feder, M., & Weiss, S. (1998). Branch Prediction Based on Universal Data Compression Algorithms. *In Proc. Twenty-Fifth International Symposium on Computer Architecture*. New York, NY: ACM.
- Fisher, J., & Freudenberger, S. (1992). Predicting Conditional Branch Directions From Previous Runs of a Program. *In Proc. Fifth Intl. Conf. on Architectural Support for Prog. Lang. and Operating Systems*. New York, NY: ACM.
- Hookway, R., & Herdeg, M. (1997). Digital FX!32: Combining Emulation and Binary Translation. *Digital Technical Journal* 9(1):3-12.
- Kistler, T., & Franz, M. (1998). Computing the Similarity of Profiling Data. *In Proc. Workshop on Profile and Feedback-Directed Optimization*.
- Wang, Z., & Rubin, N. (1998). Evaluating the Importance of User-Specific Profiling. *In Proc. 2nd USENIX Windows NT Symposium*. Berkeley, CA: USENIX Association.
- Young, C. (1997). *Path-based compilation*. Ph.D. thesis, Division of Engineering and Applied Sciences, Harvard University.