

Performance Debugging for Distributed Systems of Black Boxes

Presented by ^{Tada Yoshiakatsu} 多田 好克

SOSP 読み会 at 電気通信大学 on Nov. 15, 2003.

0 論文の概要

分散システムを ^{black-box} 暗箱として捉え、^{bottleneck} 処理の隘路発見の手助けをする方法論の提案、実践、評価。11章からなり、以下のような構成になっている。

1. 導入
2. 問題点と目標、非目標
3. 関連研究
4. アプローチの概要
5. 提案するアルゴリズム (2種類)
6. トレースの取り方
7. 実験とその結果
8. -11. 今後、サマリー、謝辞、参考文献

ヘテロで手を加えられない大きなシステムのもり。

8. ~ 11. 章は省略 ...

nesting algorithm (入れ子) と convolution algorithm (畳み込み)

1 導入

Web ベースのシステムなどの、^{multi-tired} 分散マルチタイヤシステムの虫取りや性能解析用のデータを提供。特に、様々なソフトの寄せ集めで構成され、手の出しようのない暗箱の集合体も対象にする。つまり、ソフトやネットワークに手を入れずに有用な情報を提供する。

目指すシステムはアプリ独立で、^{passive tracing} ノード間の通信を受動的に収集する。また、そのデータをオフラインで解析し、見やすく提示する。知らなくて良いし、対象とする分散システムの構造とか機能、メッセージの内容などは変更する必要もない。

って、様々なサーバを呼び合うようなシステムのこと？

もちろん、知っている、より良い情報が提供できる。

2 問題点と目標、非目標

分散システムを構成要素のサーバやクライアントと、それらを結ぶ^{delay, latency} 通信路のグラフとして捉える [図 1] ノードとエッジは、それぞれに処理時間が異なり、また、呼ばれ方によっても処理時間は異なる。ここでは、エッジの処理時間は考えない。図 1 は RPC を元にしたシステムであるが、メッセージ伝達型のシステムにも対応できる。

a) ^{causal path} 影響の大きい因果経路パターンを見つける、 b) 遅延の原因となるノードを同定する、ためのツールを提供することが目的である。

- | | |
|------------------|----------------------|
| 対象システムの知識を必要としない | × データを提供するだけ、判断は人間 |
| 対象システムに手を加えない | × 妥当なシステムを保証するものではない |
| 対象システムに外乱を与えない | × ベンチマークとしては使えない |

ツールの使い方が簡単
データの正確さは他の解析に照らして示すが、有用かどうかは読者の判断に委ねる。

ノードと呼ぶ。計算機と違って良い。

エッジと呼ぶ。要するにネットワークのこと ...

もちろん、同じように呼ばれても時間の分散はある。

LAN だから ... 仮想遅延ノードを考へても OK だし ...

e-mail システムなど ...

実行回数が多く、処理時間が長い ...

キャッシュする認証サーバは原因ではない ...

人間の判断が必要なので、示すのは難しい ...

3 関連研究

同じ方法、同じ問題意識: 文献 [11] は特別な言語を使う。magpie[13] はより複雑。また、稀な事象に重きを置いている。DPM[18] はカーネルに手を入れている。商用のシステムは JVM とか .net を対象にし、範囲が狭い。

legacy ノードには使えない。

異なる方法、同じ問題意識: 文献 [26] は実時間が目的で、プログラマが対象システムにデータ取得ルーチンを埋め込む。ETE[10] も同じで暗箱アプローチではない。

同じ方法、異なる問題意識: Pinpoint [5,6] は異常^{fault node}ノードを見つけるのが目的。文献 [3] も同様で、対象システムを止めることもある。文献 [2] は異常^{fault injection}の埋込みをするので、暗箱ではないし、対象に対するそれなりの知識が必要。文献 [27] は侵入検出が対象、文献 [12] は WAN のネットワーク解析が対象。

つまり、対象システムに外乱を与える。

4 アプローチの概要

図 2 上のような状況で、図 2 下のような情報を提供する。ツールは奇数矢印の時刻から、偶数矢印の時間を推測する。上のような呼出し関係を知らずに、2, 6, 10 の時間を提供するのが目的。そのために、以下の手順を踏む。

時刻同期は OK と思っている。また、時刻順のソートも対象外 ...
タイムアウトや再送なども含まれる。

トレース情報収集: これだけは実時間で実行する。あちこちで集めたものを時刻でソートしてまとめる。詳細は 6 章で ...

因果経路やパターン推測: ここからはオフラインの処理。トレースは大量で不完全であることが前提。処理時間、使用メモリ量は少なく、堅牢である。アルゴリズムの詳細は 5 章で ...

視覚化: あまり重要とは考えていない(手が回っていない?)

5 提案するアルゴリズム

アルゴリズムは 2 つある。1 つは RPC を対象にした入れ子^{nesting algorithm}アルゴリズム、もう 1 つはメッセージ伝達を対象にした畳込み^{convolution algorithm}アルゴリズム。

5.1 入れ子アルゴリズム

トレース情報は基本的には(時刻、送り手、受け手、call/return)の 4 つ。で、以下の 4 つの手順を踏む。

他の情報も取れれば組み入れることはできる。

呼出しペアと入れ子の同定(FindCallPairs): B C の call は $T_{opencalls}$ というハッシュに覚えて、C B の return が来れば対応づける。複数の call が先行した時には、最も遅い物から順に対応付ける。call/return がマッチしないものは捨てる。そのペアを $T_{callpairs}$ に記録。そのペアに、 $T_{opencalls}$ 内で、この時刻以前に B を call した * B の集合を記録する。

括弧内は図 4 の関数名。

つまり、入れ子にする。多くの場合、これで OK と言っている ...
問題はない ...

つまり、全ての親ネ。

(A, B, i, j) は A B が時刻 i に、B A が時刻 j に観測されたペアを表す。

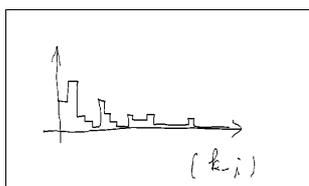
(k-i) を短い時間間隔で離散化したものを横軸にする。しかも、対数軸でやるようだ ... 1m 秒 ~ 2 時間を $\log_{1.05}$ で 340 に分けたい。

と言いながら、実践では 1 つ目だけで、後の 2 つは無視。

妥当な入れ子の決定 (ScoreNestings): ペア (B, C, k, l) は、多数のペア (A, B, i, j) に含まれる。その中から因果の深いものを同定する。方法は全ての親に対して、(k-i) を横軸にしたヒストグラム(貢献度 ...)を作る。ただし、ヒストグラム全体の面積は 1 に正規化する。図 5 の例では、4 つの場合が出るが、 $(t_3 - t_1)$ と $(t_4 - t_2)$ が同じで、他の 2 つより高い値を示す。

親の決定 (FindNestedPairs): ヒストグラムに基づいて親を決める。基本的には値の高いものを親とするが、以下の 3 種の状況では値を $1/x$ にするというペナルティを課す。

2 乗という値は経験値。



- 図 5 のように同じ親に複数回出てきた子は、その回数の 2 乗を x とする。
- 既に対応付いている親子。気持ちは、その回数の n 乗を x にしたいが、 n は 0 で x は 1、ペナルティなし。
- 親に対応付いている子の数。これも、上と同じでペナルティなし。

なお、値が同じ場合は時間の短い親に対応付ける。

呼出し経路の決定 (FindCallPath, CreatePathNode): 個々のペアに対して、親が居なければルートに登録して、CreatePathNode を呼び出す。CreatePathNode は、再帰的に実行され、ノードの実行時間と呼出しまでの時間を収集して回る。

このアルゴリズムは、(トレースしたメッセージの個数) × (ノードごとの並列度の平均) 程度の時間、空間で実行できる。

5.2 畳込みアルゴリズム

このアルゴリズムは、メッセージの集合から因果関係を見つけ出す。トレースをエッジごとのメッセージに分け、それを時刻シグナルのように処理する。出力は A B A C のような制御の因果とその処理時間。トレース情報は (時刻、送り手、受け手)。アルゴリズムは以下の手順を踏む。

1 段目の因果関係 (FindPathsFromRoot): 指定された根が送り手となるトレースの集合を求める (A *)

2 段目の因果関係 (ProcessNode): 個々の受け手 (*) について、さらに次の受け手の集合を求める (B *')。この A B *' に対して ...

畳込み (FindCausedMessages, FindCorrelation): A B と B *' の 2 つの時刻の差ごとの共起頻度をとる。つまり、時刻を短い間隔で離散化し、そこにメッセージが存在すれば 1、なければ 0 という関数 $s(t)$ を作る。で、A B と B *' の時刻を $s(t)$ で畳込むと、2 つのメッセージの時間を横軸にした共起頻度グラフ [図 7] ができる。

後は適当 ... : 図 7 のようなデータから因果の候補を決める。

- 平均より標準偏差の 4 倍以上値の大きい山を求める。
- ただし、平均より標準偏差の 3 倍以下の谷が無いと候補とはしない。

全体の処理は ... : という処理を深さ優先で再帰的にいき、因果経路を求める。

問題点として、以下のようなことが述べられている。

遅延の分散が小さいと良い結果が得られる。分散が大きい場合は図 6 の 22 行の v を大きくする。遅延の長いものと異常に短いものを除くとより良い結果が得られるが、紙面の都合で省略。また、他にも多くの機能を加えて結果を改善しているが、ここでは述べない。でも、7 章の結果はこれらの効果を反映させてある。

必要な空間は $O(m + S)$ 。 m は ProcessNode の扱うメッセージ数。 S は T / μ で、 T は最も長い因果の時間、 μ は離散化した単位時間。必要な時間は $(em + eS \log S)$ で、 e は出力のエッジ数。第 2 項の $(S \log S)$ は畳込みの FFT の計算時間。経験的には第 2 項の時間が支配的。

5.3 2 つのアルゴリズムの比較

入れ子は RPC が前提。畳込みはどのような状況にも使える。入れ子では A B C (B) A とか、遅延書き戻しのようなリターンした後に処理を依頼するものには対応できない。他方、畳込みは A B C B A のパターンに対して、A B A も出力してしまう。畳込みは稀に発生するパターンや分散の大きいパターンは検出できない。

5.4 可視化

文献 [8] の “dot” プログラムというのがあるらしい。図 8 が入れ子の出力例。子供も含むノードの実行時間、子供を呼び出すまでの平均処理時間が書かれている。また、図 9 が畳込みの出力例。メッセージの総数と、遅延時間が頻度の多い順に並べられている。

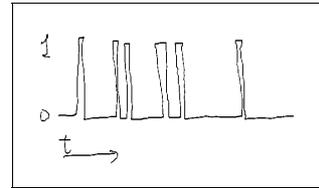
6 トレースの取り方

個々のメッセージに対するデータ量は少ないが、長時間に渡り頻度も高い。ネットワークの帯域とかスケラビリティとかが問題となる。データ収集はネットワークの受動的なモニタリング

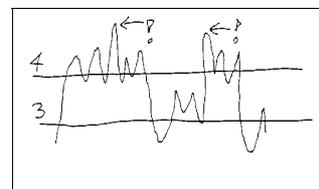
と言っても、時刻によるソートはどうするの？

パルスのトリガ関係を調べる方法と同じ ...

図 6 だけど、再帰でややこしいので、説明はプログラムに即していない。



$$(f \otimes g)_t = \sum_{x=-\infty}^{+\infty} f_x g_{t-x}$$



4 は経験的な値 ...

3 も経験的な値 ...

何！

許せん !!

つまり、入れ子を飛ばしてリターンする。

A B A の後に、B C B が実行される。

packet sniffing とも言う。

が基本だが、アプリケーション、ミドルウェア、カーネルなどに手を入れる方法でも、本論文の解析手法は有効である。

つまらん仕事ぢゃ。

port mirroring という機能を持ったスイッチがある。また、個々のホストでデータを得るという方法もある。

ネットワークの受動的なモニタリングでは、複数のパケットからメッセージを再構成する必要があるかも知れない。また、近年のスイッチではモニタリングできない。

情報源が異なるものは時刻に従って統合する。2つのアルゴリズムは時刻のズレに対してもある程度対応できる。また、分散時刻問題はNTPで何とかできると考える。

7 実験とその結果

実験は実際のトレースデータではやっていない... また、トレースは暗箱アプローチで取ったものではない。以下の様々な方法でトレースを用意した。

けど、より正確な因果経路導出と我々の方法とを比較できる。

トレース生成ソフト: maketrace コマンドを作った。様々なパターンを生成、実験できる。

J2EE: J2EE で 130 万メッセージ、3 時間のトレースを 2 つ作った。

電子メールの Received ヘッダ: 良いトレースではないが、RPC では無いものの例として ... 2 カ月で 11,683 通、81,044 の Received ヘッダを集めた。

紙面の都合と、大した結果は得られなかったからだそう。

その他: 分散ファイルシステムの物なども集めたが、以下では述べない。

以下、それぞれの結果。

7.2 マルチタイヤシステム: 図 10 のようなシステムを想定し、maketrace した。WS2 と AUTH の間に 200 m秒の遅延を入れたものとの 2 つで実験。図 11 が出力。左から頻度の高い順に並んでいる。図 12 が入れ子の、図 13 が畳みみの出力の一部 ...

だが、見るものには無い ...

7.3 J2EE: J2EE の PetStore system。 /mylist.jsp ノードの呼出しに 50 m秒の遅延を入れたものとの 2 つで実験。出力された中で重要そうな因果経路の 1 つが図 14 と図 15。

J2EE の使い方のデモソフト?

7.4 電子メールのヘッダ: 畳みみを実験。30 秒のスライスだとすべての遅延が 0 となった。スライスを 5 秒にすると、0 以外の遅延も観測された。図 16 が結果。

図 17 はマルチタイヤシステムの出現頻度の高い N パターンに対する false negative の割合。ほとんど、1/N 以下の割合に納まっている。図 18 は病的な例で、入れ子がどう間違えるかを示したものの。図 19 はメッセージの並列度が高い時の時間見積もりを失敗する様子を示したものの。図 20 は呼出し時間の分散を増加させた時の様子。図 21 はトレース生成が追いつかない時にメッセージを捨てることを想定した時の様子。1%以下ならば大丈夫。図 22 は時刻にズレがある時の様子。マルチタイヤシステムの WS2 の時刻を数 10 m秒進め、それに対して入れ子の比較に幅を持たせた (comp) と貢献度のグラフを平滑化した (smooth) 場合の様子。

詳細は省略。

表 1 は様々なメッセージ数とトレース時間による空間、時間使用量の実測値。入れ子は 1.7 G の Pen4 で Linux2.4.20 の値。畳みみは 667 M の AlphaServer で、Tru64 Unix V5.1 の値。特に畳みみの処理時間はスライスの幅に依存するが、それが $(S \log S)$ の線に乗ることを図 23 に示した。

8 ~ 11 今後、サマリー、謝辞、参考文献

8 ~ 11 章は以下の通り。

- より実践的なトレースで実験したいが、データの入手が困難。
- ロックのあるというか同期するシステムには適応できない。
- 入れ子の処理を、時間の窓を作ってその中で処理をする。処理にフェーズのあるシステムでは、短期間にだけ見られるパターンが抽出できる。また、使用空間量が減らせる。

と言いながら、単に今後の話だけが ...