
Data Recovery for Web Applications

Istemi Ekin Akkus, Ashvin Goel, *University of Toronto*
2010 IEEE/IFIP International Conference on Dependable Systems & Networks (DSN)

Sistemas Distribuídos e Tolerância a Falhas

João Pereira m3873 | Tiago Simões m3965

Contexto

- Actualmente, as aplicações Web guardam a informação no lado do servidor.
- Os *backups* podem resolver perdas de dados, mas não diagnosticam o evento causador.
- O objectivo é propor um sistema de recuperação que ajude os administradores na recuperação de dados e identificação dos eventos responsáveis.
- *“Our results show that our system enables recovery from data corruption without loss of critical data and incurs small runtime overhead.”*

Introdução

- A construção de aplicações Web é cada vez mais aberta, possibilitando integrações de *plugins* de terceiros.
- Conseqüentemente, a preocupação com o controlo de tais extensões é cada vez maior.
- A perda de dados pode ser revertida com os backups, mas as informações posteriores aos mesmos têm de ser introduzidas manualmente.

Introdução

- O sistema de recuperação proposto não depende da aplicação web.
- Assim é independente e imune às falhas ocorridas na mesma.
- Os principais objectivos do sistema são:
 1. Permitir que os administradores diagnostiquem as falhas que corromperam os dados;
 2. Permitir a recuperação selectiva dos dados sem comprometer o resto da aplicação.

Introdução

- A principal contribuição deste trabalho é a análise de dependências na recuperação de sistemas par aplicações web.
- Foi implementado um protótipo em PHP e MySQL, testado em sistemas como o Wordpress e Drupal.

Modelo da Aplicação

- **Uma aplicação web normalmente consiste em três camadas:**
 1. Apresentação;
 2. Aplicação Lógica;
 3. Base de Dados.
- **O sistema proposto tira partido das seguintes capacidades das aplicações web:**
 - Informação persistente armazenada na BD para permitir acesso concorrente, produzindo logs de acesso aos mesmos;
 - As linguagens de programação têm mecanismos próprios de controlo de excepções, etc;
 - Os web servers tratam cada pedido de forma independente, criando um processo independente, criando assim isolamento.

Considerações Prévias

- **O sistema proposto assume que:**
 - O SGBD e a linguagem de programação usada não contém bugs;
 - A corrupção dos dados está na camada da BD, resultando de bugs nas restantes camadas;
 - Se a aplicação não usa transacções, cada *query* é processada separadamente;
 - As transacções da BD podem ser replicadas correctamente, através de um nível de isolamento.

Visão Geral do Sistema

- **O sistema consiste em:**
 - Componente de monitorização on-line;
 - Dois componentes de análise e recuperação de dados após a corrupção ser identificada.
- Os componentes de análise e recuperação são usados após identificada uma corrupção.
- Os mesmos usam dados recolhidos na fase de monitorização.

Componente de Monitorização

- Responsável por correlacionar *requests* ao longo das camadas da aplicação;
- De forma resumida, uma *request* é transformada numa transação durante a recuperação;
- Esta transformação é baseada na análise do *log* permitindo o mapeamento entre *requests* e transacções da BD.

Componente de Análise

- Antes da análise, o estado actual dos dados é guardado;
- Usa dados recolhidos durante a monitorização, derivando **três tipos de dependência de dados** nas diversas camadas:
 1. Dependências da BD;
 2. Dependências da Aplicação;
 3. Dependências do Cliente.

Componente de Análise

Dependências da BD

- Ajudam na correlação das *requests* quanto às operações baseadas na BD;

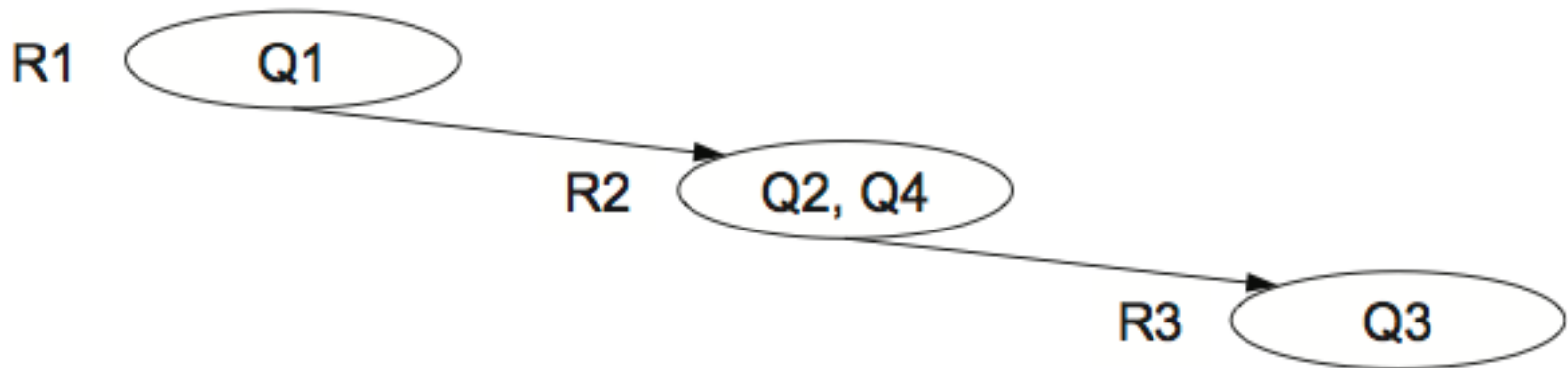


Figure 2. A request dependency graph.

Componente de Análise

Dependências da Aplicação

- Existe a possibilidade de serem geradas falsas dependências, sendo que uma *request* pode envolver múltiplas *queries* que podem não ter nenhuma dependência;
- É usada uma técnica chamada *dynamic tainting* para determinar dependências da camada lógica de aplicação relacionadas com uma determinada *request*, eliminando assim falsas dependências.

Componente de Análise

Dependências do Cliente

- Por fim, o componente de análise usa dependências do lado do cliente através das *requests*
 - Por exemplo, os *cookies* de sessão identificam todas as *requests* associadas com uma sessão de login.
- Providencia por exemplo, a possibilidade de um administrador reverter todas as acções tomadas numa dada sessão.

Componente de Recuperação

- Objectivo é disponibilizar ferramentas que facilitem a recuperação.
- Providenciam informação que permitem a identificação das *requests* responsáveis pelas falhas.
- Depois da acção de análise, este componente usa a informação do *log* da BD e da *request* para gerar transacções compensatórias.

Implementação

- Foi implementado um protótipo com PHP e MySQL.
- A maioria do código diz respeito ao componente de recuperação.
- As alterações efectuadas à estrutura do PHP e MySQL são residuais, facilitando o *port* para outras tecnologias também.
- A principal alteração ao MySQL foi a modificação do JSQParser.

Avaliação dos Resultados

Políticas de Dependência

- A avaliação de resultados tem por base a comparação da eficácia de recuperação segundo as seguintes políticas:
 1. **Request-Row:** Assume que uma *request* é *tainted* se lê uma linha *tainted* da BD.
 2. **Program-Row:** Durante uma *request* todas as variáveis que usam linhas *tainted* são marcadas como *tainted*.
 3. **Database-Row:** Propaga *taints* quando uma *query* lê uma linha *tainted* e actualiza outras linhas.
 4. **Program-Field:** Similar a 2. mas *taints* são armazenadas na BD numa granularidade do campo.
 5. **Database-Field:** Similar a 3. mas *taints* são armazenadas na BD numa granularidade do campo.

Avaliação dos Resultados

Exactidão da Recuperação

- A exactidão da recuperação, com base nas políticas definidas, é medida através da activação de 5 *bugs* reais em aplicações web, como o Wordpress, Drupal e Gallery2.
- A escolha dos *bugs* foi feita quando no mesmo:
 - Os dados foram corrompidos e só podem ser reparados com restauro do backup;
 - Os *bugs* estão relacionados com a camada de aplicação.

Avaliação dos Resultados

Exactidão da Recuperação

- Para avaliação, foram definidas acções de recuperação correctas, como sendo as acções que resolverão a corrupção dos dados.
- **Para tal foram definidas 3 métricas:**
 - Determinação das inconsistências aplicacionais responsáveis pela falha;
 - Medição de falsos positivos (marcados como responsáveis mas sem culpa);
 - Medição de falsos negativos (não marcados como responsáveis quando o deveriam ser).

Avaliação dos Resultados

Exactidão da Recuperação

Table 2. Recovery accuracy for request-level and program-level dependency policies. The false positives column shows numbers without and with table whitelisting, respectively.

Case	Total Number of Requests	Requests to Undo	Dep. Policy	False Positives	False Negatives
Wordpress - link category rename	109	1	none	0	0
			request-row	60	0
			program-row	8	0
			program-field	6	0
Drupal - lost voting information	118	7	none	0	6
			request-row	111/100	0
			program-row	95/89	0
			program-field	89/0	0
Drupal - lost comments	117	1	none	0	0
			request-row	116/102	0
			program-row	100/93	0
			program-field	95/0	0
Gallery2 - removing permissions	91	1	none	0	0
			request-row	90/13	0
			program-row	88/11	0
			program-field	82/10	0
Gallery2 - resizing images	151	1	none	0	0
			request-row	148/0	0
			program-row	139/0	0
			program-field	119/0	0

Avaliação dos Resultados

Exactidão da Recuperação

Table 3. Recovery accuracy of database-level dependency policies. All numbers indicate queries.

Case	Queries to Undo	Dep. Policy	False Positives	False Negatives	Inconsistencies after Undo
Wordpress - link category rename	23	database-row	0	15	The count value does not match the actual number of links.
		database-field	0	21	
Drupal - lost voting information	38	database-row	86	16	The poll_votes table has duplicate entries.
		database-field	0	18	
Drupal - lost comments	24	database-row	116	0	none
		database-field	0	0	
Gallery2 - removing permissions	9	database-row	97	0	The global sequence id has an old value breaking future inserts requiring a new id.
		database-field	9	0	
Gallery2 - resizing images	17	database-row	110	0	
		database-field	20	0	

Avaliação dos Resultados

Conclusões

- Estes resultados mostram que, as políticas ao nível da BD tendem a ter menos falsos positivos que os outros níveis:
 - Assim um administrador pode comparar os *outputs* da BD e determinar as melhores acções correctivas a tomar.
- Um dos maiores desafios é determinar correctamente as dependências entre *requests*;
- O prototipo implementado permite mostrar que a funcionalidade de recuperação de dados pode ser obtida com pouca sobrecarga e pouca modificação das aplicações.