

Estudio de algoritmos de calendarización para flujos de trabajo

Fernando Aguilar Reyes

Tesina para obtener el título de Ingeniero en Computación
Instituto Tecnológico Autónomo de México
Río Hondo #1, Progreso Tizapán, Del. Álvaro Obregón, 01080
México, Distrito Federal

15 de enero de 2014

Índice general

1. Introducción	1
2. Flujos de trabajo	3
2.1. Definición y ejemplos	3
2.1.1. Anotación de proteínas	3
2.1.2. Curvas de amenaza sísmica	3
2.1.3. Generación de facturas	5
2.2. Estructura de los flujos de trabajo	5
2.2.1. Perspectivas de los flujos de trabajo	6
2.3. Calendarización como control de flujo	6
3. Cómputo distribuido	9
3.1. Enfoques de cómputo para flujos de trabajo	9
3.1.1. Clusters	9
3.1.2. Grids	10
3.1.3. Nubes	10
3.2. Comparación de los enfoques	12
4. Calendarización	13
4.1. Definición del problema	13
4.2. La complejidad de calendarizar	14
4.3. Calendarización de flujos de trabajo	14
4.3.1. Reducción de flujos de trabajo a grafos dirigidos acíclicos	14
4.3.2. Definición del problema de calendarización de flujos de trabajo	14
4.3.3. Complejidad computacional de la calendarización de flujos de trabajo	15
5. Algoritmos de calendarización de flujos de trabajo	17
5.1. Algoritmos de Mejor Esfuerzo	18
5.1.1. Myopic (miope)	18
5.1.2. Min-min	18
5.1.3. Max-min	18
5.1.4. Sufragio	18
5.2. Algoritmos de Calidad en el Servicio	19

5.3.	Problemas adicionales	19
5.3.1.	Estimación del tiempo	19
5.3.2.	Descomposición de tareas	19
6.	Software para la administración y ejecución de flujos de trabajo	21
6.1.	Software orientado a clusters	21
6.1.1.	Open Grid Scheduler	21
6.2.	Software orientado a grids	21
6.2.1.	SwinDew-G	21
6.2.2.	Pegasus	21
6.3.	Software orientado a nubes	21
6.3.1.	ANEKA	21
6.3.2.	SwinDew-C	21
6.3.3.	ASKALON	21
7.	Conclusiones	23
A.	Pseudocódigos	25
A.1.	Algoritmos de menor esfuerzo	25
A.1.1.	Miope	25
A.1.2.	Definiciones para los algoritmos Max-Min y Min-min	25
A.1.3.	Min-Min	26
A.1.4.	Max-Min	26
A.1.5.	Sufragio	27
A.1.6.	HEFT	27
A.1.7.	Híbrido	28
A.1.8.	TANH	28
A.2.	Algoritmos de Calidad en el Servicio	29
A.2.1.	BackTracking	29
A.2.2.	Deadline Distribution	29
A.2.3.	CTC	29
A.2.4.	Deadline-MDP	29
A.2.5.	LOSS/GAIN	29
A.3.	Metaheurísticos	29
A.3.1.	Algoritmo Genético	29
A.3.2.	GRASP	29
A.3.3.	Recocido simulado	30

Índice de figuras

2.1.	Flujo de trabajo para la anotación de proteínas, diseñado en la Escuela Imperial de Londres para el proyecto <i>e-Protein</i>	4
2.2.	Flujo de trabajo para generar curvas de amenaza, elaborado por el SCEC. . .	4
2.3.	Flujo de trabajo para generar facturas electrónicas, utilizado por una empresa de cementos.	5
3.1.	Topología del cluster Aitzaloo	10
3.2.	Detalle de la red de fibra óptica para conectar los nodos robustos del grid del proyecto LANCAD	11

Capítulo 1

Introducción

Cada día, los cómputos son más complejos. Hay bancos que procesan millones de transacciones diarias. Las películas de animación requieren vastos tiempos de cómputo. Muchos proyectos de computación científica requieren hacer numerosos cálculos para llegar a resultados pertinentes. ¿Qué tienen en común todas estas aplicaciones? Ya hemos mencionado que toman mucho tiempo calcularse. Entonces, es natural pensar que se puede distribuir el gran trabajo que requieren estos proyectos entre varias computadoras para disminuir el tiempo total de ejecución y de esta forma, tener una solución escalable.

Lograr esta paralelización requiere un esfuerzo por parte del desarrollador de la solución. Aunque estas técnicas de paralelización son muy efectivas, éstas son aplicadas cuando el problema a resolver ha sido bien definido y cuando sólo hay una instancia del problema. Ahora bien, cuando la solución involucra varios pasos que están relacionados entre sí, por ejemplo, que el programa C requiera de la salida del programa A y del programa B para que pueda funcionar. O, cuando estamos definiendo los grandes bloques de solución del problema, hay una cierta secuencia que debemos seguir y, dentro de dicha secuencia, hay algunos pasos que podemos resolver de manera concurrente. Por lo tanto, estamos haciendo un *modelo* de nuestro problema.

Este modelo también es llamado *flujo de trabajo*. De manera muy abstracta, un modelo de trabajo es un conjunto de pasos que modelan la ejecución de un proceso. Este concepto ha sido aplicado en numerosas áreas. En el ámbito de los negocios, se utiliza diagramas dibujados con los bloques y reglas del *Business Process Execution Language* para modelar procesos de negocio.

En el ámbito de las ciencias en computación, los flujos de trabajo son aplicados para modelar problemas que requieran varios pasos para su solución. Algunos flujos de trabajo requieren mucho tiempo de ejecución para sus pasos, como son el caso de los flujos de trabajo científicos. Otros flujos de trabajo son muy simples pero se requieren que se ejecuten muchos de éstos. De esta forma, es deseable distribuir la ejecución de éstos flujos de trabajo entre varias computadoras. Si bien es posible paralelizar algunos pasos de la ejecución de nuestro flujo de trabajo utilizando las técnicas antes mencionadas, hay restricciones de orden que se deben respetar, por lo cual, se debe planear la ejecución del flujo de trabajo.

Hay varias maneras de hacer esta planeación de la ejecución, también llamada *calendari-*

zación. Cabe aclarar que definimos la calendarización como la asignación de recursos a tareas de tal modo que se ejecuten todas las tareas. Con ello, se desea encontrar una forma óptima de hacer esta calendarización, como reducir el tiempo de ejecución total del flujo de trabajo. Sin embargo, con la aparición del cómputo en la nube, es posible ejecutar nuestro flujo de trabajo con otras restricciones, como minimizar el presupuesto necesario para la ejecución del flujo afectando el tiempo de ejecución.

En esta tesina se hace un estudio de los principales algoritmos de calendarización de flujos de trabajo, con énfasis en los algoritmos utilizados en cómputo distribuido, en especial en cómputo en la nube. En el capítulo 2 se habla de un estudio detallado del concepto de flujos de trabajo y su aplicación en computación. El capítulo 3 trata los principales enfoques de cómputo distribuido para ejecutar estos flujos. En el capítulo 4 se hace un estudio de los principales algoritmos de calendarización de los flujos. Finalmente, en el capítulo 6 discutiremos algunas conclusiones sobre el análisis de estos algoritmos.

Capítulo 2

Flujos de trabajo

2.1. Definición y ejemplos

Un flujo de trabajo es un conjunto de pasos que modelan la ejecución de un proceso [6]. En la literatura especializada, los flujos de trabajos también son conocidos como *workflows*. En particular, se estudian a los flujos de trabajo utilizados para vislumbrar la ejecución de un proceso de cómputo. A continuación, se muestran algunos ejemplos de estos flujos de trabajo.

2.1.1. Anotación de proteínas

En el proyecto *e-Protein*, realizado por la Escuela Imperial de Londres, se realizó un flujo de trabajo para la anotación de proteínas. El objetivo del proyecto [10] era la identificación y anotación de partes de proteínas que expliquen su estructura y su función. En la figura 2.1 se muestra el flujo de trabajo desarrollado, donde las cajas representan los programas que son ejecutados para cada paso del proceso de anotación, y las líneas que conectan a las cajas representan las dependencias de datos entre los programas, es decir, si una caja tiene una línea que apunta a ella, significa que dicho programa depende de otro programa determinado por el otro extremo de la flecha.

2.1.2. Curvas de amenaza sísmica

Otro notable ejemplo es el proceso para generar curvas de amenaza que describen las probabilidades de que ocurra un temblor en una determinada área. Para elaborar estas curvas, los científicos del Centro de Terremotos del Sureste de California (SCEC por sus siglas en inglés) tienen que realizar una gran cantidad de simulaciones para que sus resultados puedan ser combinados y expresados en la curva de amenaza [5]. El flujo de trabajo para generar la curva de amenaza se muestra en la figura 2.2.

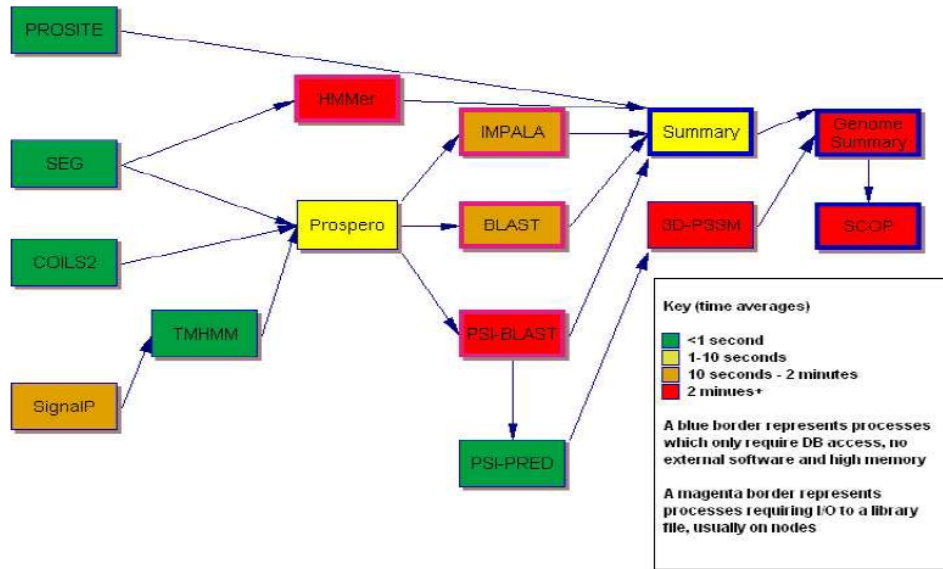


Figura 2.1: Flujo de trabajo para la anotación de proteínas, diseñado en la Escuela Imperial de Londres para el proyecto *e-Protein*.

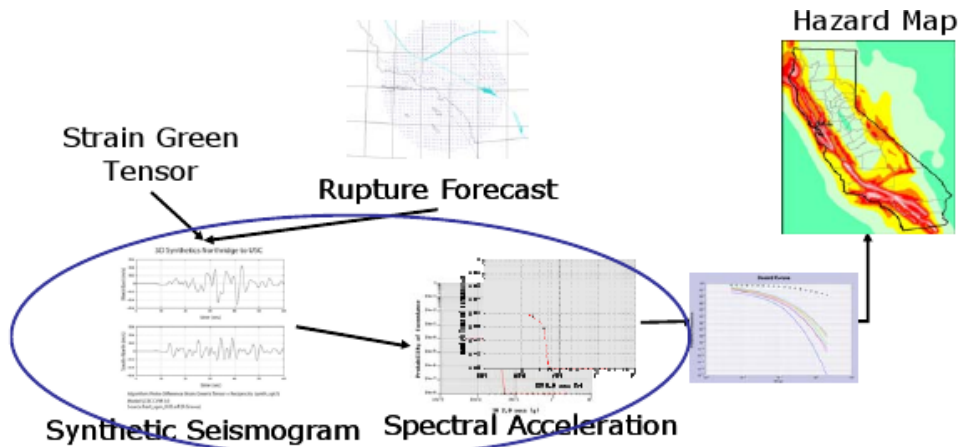


Figura 2.2: Flujo de trabajo para generar curvas de amenaza, elaborado por el SCEC.

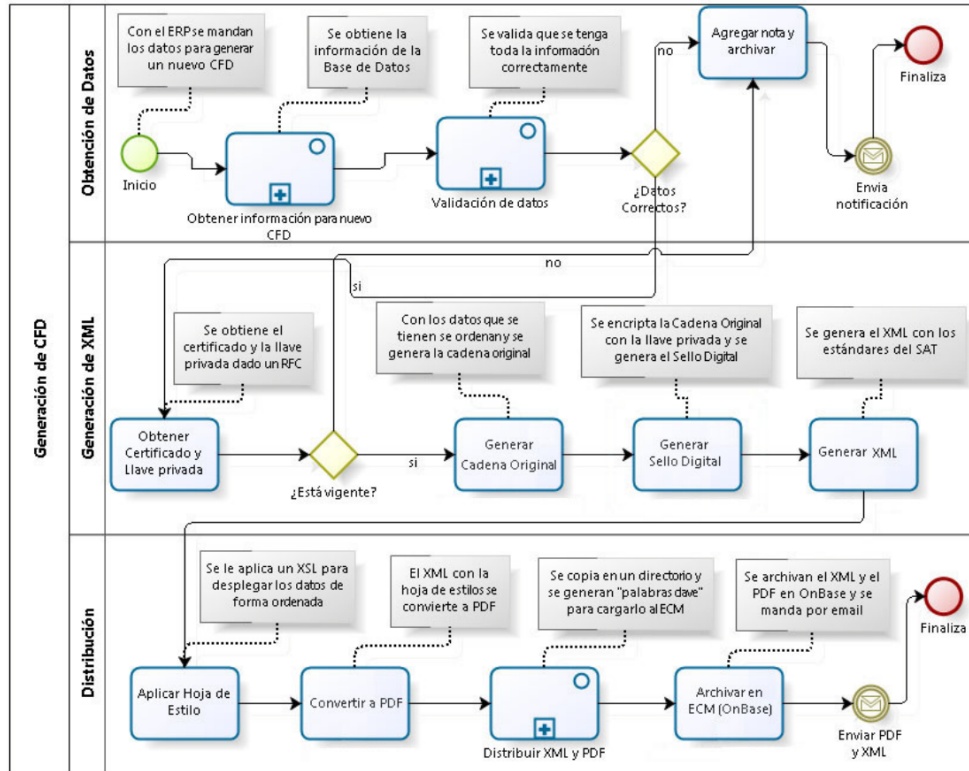


Figura 2.3: Flujo de trabajo para generar facturas electrónicas, utilizado por una empresa de cementos.

2.1.3. Generación de facturas

Recientemente en México, el Sistema de Administración Tributaria emitió los lineamientos para que las operaciones de compra-venta entre personas físicas y morales puedan ser registradas por medio de facturas electrónicas. Para generar estos Comprobantes Fiscales Digitales, las empresas tienen que hacer procesos de validaciones de RFC y encriptar el contenido de la factura con un mecanismo de llave privada. Naturalmente, este proceso de generación de factura requiere de varias actividades. En la figura 2.3 se muestra el flujo de trabajo que utiliza una empresa para generar sus facturas. Este flujo se encuentra documentado en una tesina presentada por Alcerreca [3].

2.2. Estructura de los flujos de trabajo

Como hemos visto en los dos ejemplos anteriores, los flujos de trabajo presentados describen los *pasos* que necesitan para calcular la solución de un problema. En éstos, no se ha hablado sobre detalles para hacer los cálculos necesarios para cada flujo. Tampoco definen las plataformas de cómputo en que se ejecutan estos flujos. Tan sólo definen los grandes

pasos para solucionar el problema y las dependencias que tienen estos pasos.

Ahora, es muy común que estas dependencias estén dictadas por los datos que requiere cada paso para funcionar. Sin embargo, hay situaciones en las que los pasos del flujo no requieren los datos del paso anterior para funcionar, sino que las dependencias están marcadas por el orden temporal que deben seguir estos pasos. Así, se puede notar una ambigüedad en la definición de un flujo de trabajo.

2.2.1. Perspectivas de los flujos de trabajo

El trabajo de Van der Aalst et al. [13] identifica una serie de perspectivas vislumbrados en las definiciones de las especificaciones de un flujo de trabajo en las que se basan los sistemas que administran la ejecución de éstos. Las cuatro perspectivas se enuncian a continuación:

- **Perspectiva de control de flujo.** Describe las relaciones de las actividades (pasos) con estructuras de control, tales como: secuencia, decisión, ejecución de actividades en paralelo y punto de sincronización conjunta¹.
- **Perspectiva de datos.** En ella, los flujos de trabajo describen las entradas y salidas de datos, tanto de ejecución como de control, que se tienen en cada actividad del flujo. También se toma en cuenta los datos locales a cada actividad; es decir, que sólo son necesarios dentro del contexto de ésta.
- **Perspectiva de recursos.** Muestra cuáles son los recursos con los que se cuentan para ejecutar el flujo de trabajo y la forma en que estos recursos se encuentran organizados. Estos recursos pueden ser desde entidades de cómputo hasta roles con responsabilidades específicas cumplidas por actores humanos.
- **Perspectiva operacional.** Aquí se detallan las operaciones elementales necesarias en cada actividad para ejecutar el flujo de trabajo. Estos detalles incluyen las transferencias de datos entre las operaciones y su correspondencia en programas.

Cabe aclarar que estas perspectivas están relacionadas entre sí de modo que el control de flujo es la base en la que descansan las demás perspectivas. Esto es porque la perspectiva de datos requiere que el control de flujo tenga los datos de entrada y salida como prueba de que se cumplieron las pre y post-condiciones de cada actividad, respectivamente; la perspectiva de recursos define con qué se ejecutarán y almacenará los datos del flujo de trabajo; mientras que la perspectiva operacional trata los detalles sobre cómo se utilizan *físicamente* los recursos, los datos y los programas a lo largo del flujo de trabajo.

2.3. Calendarización como control de flujo

El hecho de que la perspectiva de control de flujo sea la base de las demás perspectivas indica que la forma en que controla la ejecución del flujo determina de manera fundamental

¹Esta estructura de control también es conocida como *join synchronization*

el rendimiento de la ejecución total de una instancia del flujo. Por lo tanto, es de vital importancia encontrar métodos de calendarización que permitan encontrar correspondencias entre recursos y actividades que cumplan con los requisitos dictados en las especificaciones examinadas en la perspectiva del control del flujo y que también maximicen el rendimiento de la ejecución de todo el flujo en general.

En este trabajo, se establecerán las siguientes consideraciones para limitar el estudio de los algoritmos de calendarización, a saber:

- Los flujos de trabajo están representados como *grafos dirigidos acíclicos*, con el objetivo de simplificar el estudio
- Las tareas (o actividades) de los flujos de trabajo son consideradas *atómicas*, i.e., una tarea no puede ejecutarse incompletamente. Algunos sistemas de administración de ejecución de flujos de trabajo tienen mecanismos para lidiar con estos errores, pero el estudio de éstos queda fuera de los límites de este trabajo.
- La información de las *demás perspectivas* se puede requerir, dependiendo si el mecanismo de calendarización lo considere necesario para tomar mejores decisiones.

Capítulo 3

Cómputo distribuido

Los flujos de trabajo requieren de recursos de cómputo para su ejecución. En los ejemplos anteriores se mencionaron aplicaciones científicas que, por su naturaleza, requieren vasto tiempo de ejecución para llevarse a cabo. Por otro lado, también existen aplicaciones de negocio que, si bien son computacionalmente sencillas, se requiere ejecutar una gran cantidad de instancias de estos flujos de trabajo. Por ello, ejecutar cualquiera de estos flujos de trabajo en una sola computadora resulta prohibitivo. Así, se hace uso de varias computadoras para distribuir el esfuerzo para correr estos grandes procesos.

Dependiendo de cómo estén organizados estas computadoras, se definen los enfoques para correr los flujos de trabajo con cómputo distribuido.

3.1. Enfoques de cómputo para flujos de trabajo

Diversos enfoques se han aplicado para distribuir la ejecución de un flujo de trabajo entre varias computadoras. De acuerdo a Buyya et al., los enfoques de cómputo más importantes para los flujos de trabajo son los *clusters*, los *grids* y las *nubes* [4]. A continuación, explicaremos cada uno de los enfoques.

3.1.1. Clusters

Los *clusters* son sistemas distribuidos, paralelos, compuestos de varias computadoras que son vistas como un único recurso de cómputo [4].

Un ejemplo de un cluster es la instalación de la Universidad Autónoma Metropolitana, campus Iztapalapa, llamada *Aitzaloo*, compuesta por 270 nodos de cómputo, cada uno equipado con dos procesadores Intel Xeon Quad-Core y 16GB en RAM; los nodos están conectados entre sí por medio de switches Ethernet e Infiniband. El cluster también cuenta con un sistema de archivos distribuido basado en Lustre. La capacidad real de cómputo del cluster Aitzaloo es de 18.4 teraFLOPS [2].

El detalle de la topología del cluster se puede apreciar en la figura 3.1, en donde podemos apreciar que los switches son los puntos de conexión entre el nodo maestro, el sistema de

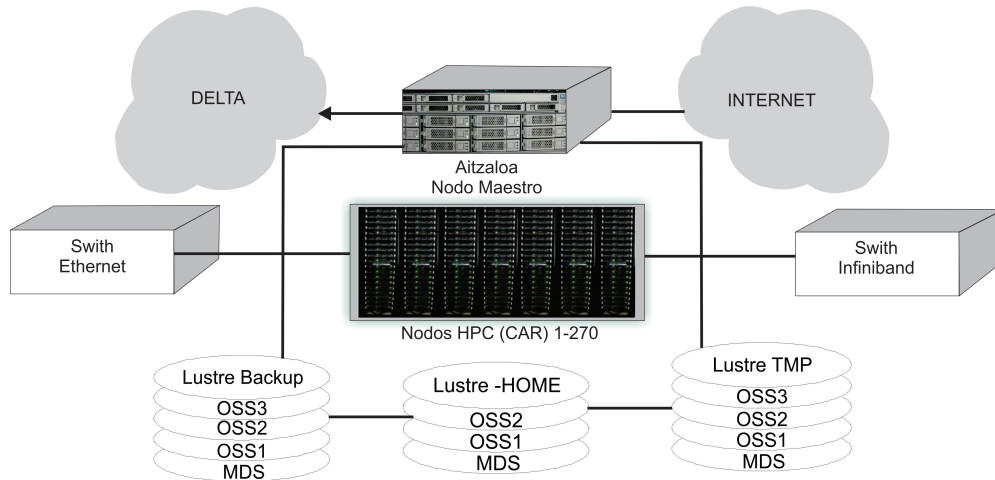


Figura 3.1: Topología del cluster Aitzaloo

almacenamiento distribuido y los nodos de cómputo.

3.1.2. Grids

Los *grids* son sistemas distribuidos, paralelos, compuestos de computadoras autónomas y geográficamente distribuidas que pueden trabajar en conjunto o de manera independiente de acuerdo a los objetivos, políticas y mecanismos de uno o varios administradores del sistema, es decir, un grid puede ser compartido entre varias instituciones [4].

El proyecto *LANCAD*¹ es un buen ejemplo, pues une el cluster *Aitzaloo* de la UAM, el cluster de la UNAM *KamBalam*, y el cluster *Xiuhcoatl* del CINVESTAV por medio de una red de fibra óptica instalada en las estaciones del Sistema de Transporte Colectivo Metro. La suma de la potencias reales de cada *nodo robusto* del grid es de 48.55 teraFLOPS [1].

En la figura 3.2 se muestra los tramos de línea de Metro que cuentan con fibra óptica para conectar cada uno de las supercomputadoras (clusters) de las tres instituciones.

3.1.3. Nubes

Las *nubes* (clouds) son sistemas distribuidos, paralelos, compuestos de computadoras o máquinas virtuales interconectadas que son aprovisionadas para usarse como uno o varios recursos de cómputo, de acuerdo a un contrato de nivel de servicio acordado entre el proveedor de la nube y el cliente [4].

Empresas nuevas y existentes proveen servicios de cómputo en la nube, tales como Go-Grid, Rackspace, Amazon, Microsoft, IBM, Oracle, entre otras. La forma en que operan es la siguiente: se paga cierta cantidad por utilizar servicios de cómputo o almacenamiento durante determinado tiempo. Así, los clientes no tienen que invertir grandes cantidades de

¹Laboratorio Nacional de Cómputo de Alto Rendimiento

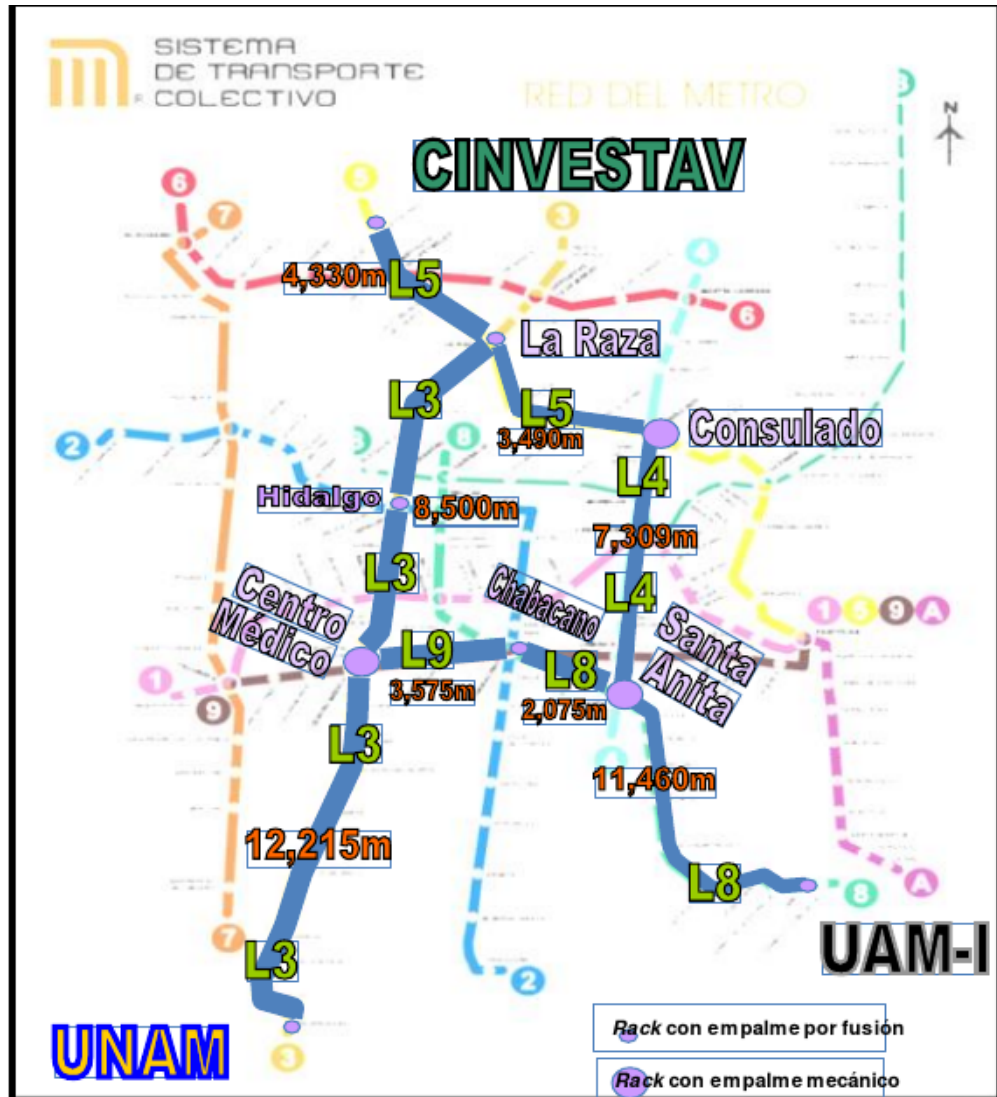


Figura 3.2: Detalle de la red de fibra óptica para conectar los nodos robustos del grid del proyecto LANCAD

dinero para contar con una gran infraestructura como en el caso de los clusters y los grids.

3.2. Comparación de los enfoques

Si bien estos enfoques difieren, principalmente, en la forma en que están organizadas las unidades de cómputo, podemos enumerar algunas observaciones.

Primero, es natural ver que los grids están compuestos de clusters, pero también podría verse a un grid como un gran cluster. Hasta cierto punto este razonamiento parece correcto. Sin embargo, la diferencia importante entre clusters y grids radica en el hecho de que el grid es comúnmente administrado por varios técnicos de varias instituciones que, pueden tener objetivos y necesidades diversas, mientras que un cluster requiere de menos administradores y se asume una mayor flexibilidad para enfrentar los fallos de un cluster.

Segundo, tanto los grids como las nubes son físicamente muy similares: ambos consisten de varios clusteres interconectados y distribuidos geográficamente. Entonces, ¿cuál es la diferencia? Lo que distingue a las nubes de los grids son dos características: (1) en una nube, se utilizan *máquinas virtuales* para distribuir el recurso. El usuario tiene la vista de una parte del grid como si fuera una computadora dedicada exclusivamente para éste. Por otro lado, en los grids, es común que el usuario acceda a él con una cuenta asignada por un administrador, de tal suerte que se tiene una vista de una gran máquina que es compartida entre varios usuarios. Muchos de los administradores de grids en el mundo administran los trabajos encargados por los usuarios de tal modo que se ejecuten lo más rápido posible con los recursos disponibles en un tiempo dado. Los usuarios, en algunos casos obtienen acceso al grid y trabajan sus proyectos de manera gratuita o, en caso contrario, pagan una cuota por un tiempo fijo de cómputo. Pero, en el caso de las nubes, se agrega un modelo económico para acceder a los recursos llamado *pay as you go*, en donde el usuario paga una cuota por los recursos utilizados, sin alguna limitación mas que el presupuesto. Así, el usuario podría pagar una gran cantidad para que su flujo de trabajo se ejecute en el menor tiempo posible o también podría pagar una menor cantidad sacrificando el tiempo total de ejecución del flujo.

Capítulo 4

Calendarización

La calendarización¹ es el proceso de asignación de recursos a tareas, de tal modo que se define un orden de ejecución de las tareas, teniendo lugar diferentes combinaciones de recursos y tareas. En este capítulo se define de la forma más elemental el problema de la calendarización para estudiar sus propiedades y también se explica la relación de este problema fundamental con los flujos de trabajo.

4.1. Definición del problema

De acuerdo a Ullman et al. [12], el problema de la calendarización se define de la siguiente manera:

Definición 1. *El problema de la calendarización consiste en:*

1. *Un conjunto de tareas $S = \{J_1, J_2, \dots, J_n\}$*
2. *Un ordenamiento parcial \prec sobre S*
3. *Una función de costo $W : S \mapsto \mathbb{Z}^+$, la cual indica el tiempo que tarda en completarse cada una de las tareas en S*
4. *Un número de computadoras (procesadores) k*

Para este caso, se asume que cuando una computadora ejecuta una tarea, ésta es ejecutada completamente. Así, el objetivo es minimizar el tiempo total de ejecución, denotado por t_{max} , respetando el orden parcial definido por \prec , asignando tareas de S a los k procesadores .

¹También conocida en la literatura especializada como *scheduling*

4.2. La complejidad de calendarizar

Como puede notarse, hay varias maneras de acomodar las k computadoras para que se ejecuten todas las tareas de S . Sin embargo, Ullman et al. han demostrado que este problema pertenece a la categoría NP-completo [12].

Esto significa que no se ha encontrado un algoritmo que pueda resolver el problema en tiempo polinomial. Entonces, la solución ingenua de probar ordenadamente todas las posibles asignaciones de tareas a computadoras resulta computacionalmente muy caro.

Así, la forma de atacar estos problemas NP-completos es utilizar métodos de aproximación [7] que obtengan soluciones subóptimas o utilizar heurísticas que resuelvan este problema sumando ciertas restricciones.

4.3. Calendarización de flujos de trabajo

Hasta ahora, se ha mencionado el problema básico de calendarización con restricciones. Con ello, se pretende plantear el problema de la calendarización de flujos de trabajo y apoyarse en la definición del problema básico de calendarización para enumerar propiedades sobre este problema.

Para plantear el problema de la calendarización de flujos de trabajo, primero se demostrará que bajo ciertas condiciones, un flujo de trabajo puede ser reducido a un grafo dirigido acíclico haciendo las transformaciones adecuadas. Luego, se utilizará la definición del problema básico de calendarización con restricciones y su semejanza para flujos de trabajo. Finalmente, se hará una descripción de la complejidad de calendarizar flujos de trabajo.

4.3.1. Reducción de flujos de trabajo a grafos dirigidos acíclicos

En el trabajo de de Mair et al. [9] proponen un formato para una representación intermedia de flujos de trabajo, con el fin de transformar una especificación de un flujo de trabajo detallada –escrita en un lenguaje basado en XML llamado AGWL– a una representación intermedia basada en grafos dirigidos acíclicos.

Aunque no existe un consenso general sobre cuál es una definición completa de un flujo de trabajo [13], el lenguaje AGWL es una especificación suficiente para una gran cantidad de flujos de trabajo de ámbito científico, porque con las construcciones if, while, for, parallel, secuencia son adecuadas para expresar vastos flujos.

El hecho importante a recalcar es que los grafos dirigidos acíclicos sirven como representación para utilizarlos en los algoritmos de calendarización.

4.3.2. Definición del problema de calendarización de flujos de trabajo

Una vez que se ha establecido a los grafos dirigidos acíclicos como nuestra representación básica de flujos de trabajo, se definirá el problema que conlleva asignar recursos a las tareas

del flujo.

Con el fin de no perder generalidad, tomaremos la definición de Wieczorek-Prodan [14] del problema:

Definición 2. *Un **flujo de trabajo** es un grafo dirigido $w \in \mathcal{W}$, $w = (\mathcal{V}, \mathcal{E})$, compuesto de un conjunto de nodos \mathcal{V} y un conjunto de aristas \mathcal{E} , donde los nodos y las aristas representan tareas $\tau \in \mathcal{T}$ y transferencias de datos $\rho \in \mathcal{D}$.*

Hay que notar que en la definición anterior que la forma en que se relacionan las tareas y las transferencias de datos con los nodos y aristas del grafo está determinada por la representación *concreta* del flujo de trabajo.

Ahora, se definirán los recursos de cómputo en donde se ejecuta el flujo. La definición de Wieczorek-Prodan habla de grids, pero es fácilmente aplicable a otros enfoques de cómputo

Definición 3. *Un **servicio** es una entidad de cómputo que puede ejecutar una tarea $\tau \in \mathcal{T}$. El conjunto de todos los servicios disponibles para ejecutar el flujo de trabajo está denotado por \mathcal{S} .*

De este modo, la calendarización es definida como una función:

Definición 4. *La **calendarización de un flujo de trabajo** w es una función $f : \mathcal{T} \mapsto \mathcal{S}$ que asigna servicios a las tareas del flujo. El conjunto que contiene todas las posibles calendarizaciones de w es denotado por \mathcal{F} .*

Cada posible calendarización determina un costo de ejecución. A continuación, definiremos el modelo de costo determinado por la utilización de los servicios.

Definición 5. *Un **modelo de costos** es un conjunto de criterios $C = \{c_1, c_2, \dots, c_n\}$ que determinan las restricciones en las que se debe ejecutar una calendarización, por ejemplo, un límite en el tiempo de ejecución, en el costo monetario o la tolerancia a fallos, entre otros.*

Definición 6. *Para cada criterio $c_i \in C$, existe una **función de costo parcial** $\Theta_i : \mathcal{S} \mapsto \mathbb{R}$, en la que a cada servicio disponible para ejecutar el flujo de trabajo se le asocia un costo por ejecutar dicho servicio con las restricciones dictadas con el criterio c_i .*

Definición 7. *Para cada criterio $c_i \in C$, existe una **función de costo total** $\Delta_i : \mathcal{W} \times \mathcal{F} \mapsto \mathbb{R}$, que asigna un flujo de trabajo w calendarizado por f un costo basado en los costos parciales determinados por los servicios utilizados para ejecutar las tareas del flujo.*

El objetivo del problema de la calendarización de flujos de trabajo es encontrar la calendarización f que minimice las funciones de costo total Δ_i , $1 \leq i \leq n$.

4.3.3. Complejidad computacional de la calendarización de flujos de trabajo

Después de haber enunciado varias definiciones para definir el problema de la calendarización de flujos de trabajo, se hará una analogía con el problema básico de calendarización para poder estudiar su complejidad computacional. La analogía es la siguiente:

1. El conjunto de tareas \mathcal{S} equivale a las tareas \mathcal{T} descritas por el flujo de trabajo.
2. El ordenamiento parcial \prec está representado tanto por las dependencias de datos \mathcal{D} y las aristas \mathcal{E} que representan el control de flujo que debe respetarse para ejecutar el flujo.
3. Las k computadoras equivalen al conjunto de servicios \mathcal{S} disponibles para ejecutar el flujo.
4. La función de costo W tiene una equivalencia implícita. El problema básico asume que conocemos apriori el tiempo de ejecución de una tarea. Sin embargo, es muy frecuente que sólo tengamos una estimación del tiempo de ejecución. Por otro lado, podemos establecer una función auxiliar que relacione las tareas con los servicios. Dicha relación es la función de calendarización f . En efecto, ejecutar una tarea en un servicio genera un costo, determinado por las funciones parciales y totales. Por lo pronto, estableceremos una relación proporcional entre costo y tiempo, con el fin de mostrar que existe una equivalencia entre la función de costo W del problema básico y el modelo de costo de un flujo de trabajo.

Con lo anterior, hemos establecido que el problema de la calendarización de flujos de trabajo es equivalente al problema básico de calendarización. Entonces, es natural inferir que el problema de la calendarización de flujos de trabajo pertenece a la categoría de problemas NP-completo.

Capítulo 5

Algoritmos de calendarización de flujos de trabajo

En el capítulo anterior definimos el problema de calendarizar flujos de trabajo. También se vio que dicho problema pertenece a la categoría de los problemas NP-completos, lo cual significa que la complejidad –o tiempo de ejecución– para resolver este problema no está acotada por una función polinomial. Por ello, diversos algoritmos se han propuesto para atacar, ya sea de manera total o parcial, el objetivo principal de la calendarización: optimizar el tiempo de ejecución total.

Ahora, con los enfoques de cómputo propuestos en el capítulo 2, se crean diversos algoritmos para diferentes necesidades. Mientras que en los clusters y los grids se asumen que estos recursos son compartidos, los algoritmos de calendarización para estos recursos asumen que los flujos deben ser ejecutados lo más pronto posible, con el fin de hacer que el recurso esté ocupado la mayor parte del tiempo. Por otro lado, el enfoque de nubes permite diseñador de flujos de trabajo elegir entre ejecutar el flujo con todos los recursos a costa de un elevado presupuesto o, minimizar dicho presupuesto tolerando un tiempo de ejecución mayor al mínimo posible.

De acuerdo a Yu et al. [15], se pueden clasificar los algoritmos de calendarización los algoritmos de flujos de trabajo ejecutados en grids en dos grandes niveles: los algoritmos de Mejor Esfuerzo y los algoritmos de Calidad en el Servicio¹. Al primer grupo pertenecen aquellos algoritmos que tratan de minimizar el tiempo total de ejecución², haciendo uso de todos los recursos disponibles. El segundo grupo, los algoritmos tratan de obtener una calendarización que cumpla las restricciones especificadas como una medida de calidad, con la posibilidad de elegir soluciones que tomen un tiempo de ejecución subóptimo.

A continuación, mostraremos los algoritmos descritos en el trabajo de Yu et al., con los ajustes en notación necesarios para que coincidan con las definiciones de flujo de trabajo establecidas en el capítulo anterior.

¹En la literatura especializada se conoce a este término como *Quality of Service*

²También conocido como *makespan*

5.1. Algoritmos de Mejor Esfuerzo

Este tipo de algoritmos tratan de minimizar algún criterio, que en muchos casos es el tiempo total de ejecución o *makespan*. Cabe aclarar que para los siguientes algoritmos, se asumirá que el grafo del flujo de trabajo tiene una correspondencia biyectiva entre los nodos *mathcal{V}* y las tareas *mathcal{T}*, es decir, cada nodo del grafo representa una tarea del flujo de trabajo. Las aristas del grafo representan *dependencias entre tareas*.

5.1.1. Myopic (miope)

Este es el más simple de todos los algoritmos. Lo único que hace es buscar un recurso disponible que pueda ejecutar la tarea y asignarle dicha tarea. No toma en cuenta otra característica a optimizar. Este algoritmo fue propuesto por Ramamritham et al. [11]. El algoritmo [15] presentado en este trabajo se encuentra en la sección [?] del Apéndice.

5.1.2. Min-min

El algoritmo min-min está basado en la heurística de terminar las tareas más cortas en el menor tiempo posible. Para ello, hace una estimación del tiempo de ejecución tomando en cuenta el tiempo de preparación de las tareas en los servicios –o recursos– para tener el tiempo y los archivos necesarios disponibles para el recurso en cuestión. El algoritmo fue presentado por Maheswaran et al. [8]. El pseudocódigo está descrito en el listado A.1.3.

5.1.3. Max-min

El algoritmo max-min –también propuesto Maheswaran et al. [8]– es muy similar al algoritmo min-min. La diferencia radica en que éste calendariza tareas cuyo tiempo mínimo de ejecución es el mayor, de tal modo que se ejecutan las tareas más *largas*.

El único cambio que se necesita hacer al algoritmo A.1.3 es cambiar la siguiente línea (14):

$$T \leftarrow \arg \min_{t \in \text{availTasks}} ECT(t, r);$$

por

$$T \leftarrow \arg \max_{t \in \text{availTasks}} ECT(t, r);$$

Así, el algoritmo modificado puede verse en el pseudocódigo A.1.4.

5.1.4. Sufragio

El algoritmo Sufragio³ [8] es una variación del algoritmo Min-min, el cual considera el valor del sufragio para hacer la calendarización. Dicho valor es la diferencia entre el menor tiempo de ejecución para una tarea *t* sobre un conjunto de recursos disponibles y el segundo

³También conocido como *Sufferage*

menor. Se calendariza a la tarea que tenga el valor del sufragio más alto, por el hecho de que las tareas que son muy sensibles a los cambios de los recursos deben ser calendarizadas primero. El algoritmo se encuentra en la sección A.1.5 del Apéndice.

5.2. Algoritmos de Calidad en el Servicio

5.3. Problemas adicionales

5.3.1. Estimación del tiempo

5.3.2. Descomposición de tareas

Capítulo 6

Software para la administración y ejecución de flujos de trabajo

6.1. Software orientado a clusters

6.1.1. Open Grid Scheduler

6.2. Software orientado a grids

6.2.1. SwinDew-G

6.2.2. Pegasus

6.3. Software orientado a nubes

6.3.1. ANEKA

6.3.2. SwinDew-C

6.3.3. ASKALON

Capítulo 7

Conclusiones

Apéndice A

Pseudocódigos

A.1. Algoritmos de menor esfuerzo

A.1.1. Miope

Input: Un grafo de flujo de trabajo $w = (\mathcal{V}, \mathcal{E})$

Output: Una calendarización f

- 1: **while** $\exists v \in \mathcal{V}$ no completada **do**
- 2: $t \leftarrow$ Obtener una tarea lista, no calendarizada, con padres calendarizados;
- 3: $r \leftarrow$ Obtener un recurso que pueda resolver la tarea en el menor tiempo;
- 4: Calendarizar t en r , i.e., $f(r) = t$;
- 5: **end while**

A.1.2. Definiciones para los algoritmos Max-Min y Min-min

Los algoritmos min-min y max-min utilizan las siguientes estimaciones:

- $EET(t, r)$ – **Tiempo estimado de ejecución:** Tiempo que el recurso (servicio) r tomará en ejecutar la tarea t , desde que la tarea es ejecutada en el recurso
- $EAT(t, r)$ – **Tiempo estimado de disponibilidad:** Tiempo en el que el recurso r estará disponible para ejecutar la tarea t
- $FAT(t, r)$ – **Tiempo de archivo disponible:** Tiempo más pronto en que todos los archivos requeridos por la tarea t están disponibles en el recurso r
- $ECT(t, r)$ – **Tiempo estimado de terminación:** Tiempo estimado en cual la tarea t terminará su ejecución en el recurso r :

$$ECT(t, r) = EET(t, r) + \max(EAT(t, r), FAT(t, r))$$

- $MCT(t)$ – **Tiempo mínimo estimado de terminación:** ECT mínimo para la tarea t sobre todos los recursos disponibles, es decir:

$$MCT(t) = \min_{r \in \mathcal{S}} ECT(t, r)$$

A.1.3. Min-Min

Input: Un grafo de flujo de trabajo $w = (\mathcal{V}, \mathcal{E})$

Output: Una calendarización f

```

1: while  $\exists v \in \mathcal{V}$  no completada do
2:    $t \leftarrow$  Obtener conjunto de tareas listas, no calendarizadas, con padres calendarizados;
3:   SCHED( $tasks$ );
4: end while
5: procedure SCHED( $availTasks$ )
6:   while  $\exists t \in availTasks$  no calendarizadas do
7:     for  $t \in availTasks$  do
8:        $res \leftarrow$  Obtener recursos disponibles para  $t$ ;
9:       for  $r \in res$  do
10:        Calcular  $ECT(t, r)$ ;
11:      end for
12:       $R_T \leftarrow \arg \min_{r \in res} ECT(t, r)$ ;
13:    end for
14:     $T \leftarrow \arg \min_{t \in availTasks} ECT(t, R_T)$ ;
15:    Calendarizar  $T$  en  $R_T$ ;
16:    Remover  $T$  de  $availTasks$ ;
17:    Actualizar  $EAT(R_T)$ ;
18:   end while
19: end procedure
```

A.1.4. Max-Min

Input: Un grafo de flujo de trabajo $w = (\mathcal{V}, \mathcal{E})$

Output: Una calendarización f

```

1: while  $\exists v \in \mathcal{V}$  no completada do
2:    $t \leftarrow$  Obtener conjunto de tareas listas, no calendarizadas, con padres calendarizados;
3:   SCHED( $tasks$ );
4: end while
5: procedure SCHED( $availTasks$ )
6:   while  $\exists t \in availTasks$  no calendarizadas do
7:     for  $t \in availTasks$  do
8:        $res \leftarrow$  Obtener recursos disponibles para  $t$ ;
9:       for  $r \in res$  do
10:        Calcular  $ECT(t, r)$ ;

```



```

11:         end for
12:          $R_T \leftarrow \arg \min_{r \in res} ECT(t, r);$ 
13:     end for
14:      $T \leftarrow \arg \max_{t \in availTasks} ECT(t, R_T);$ 
15:     Calendarizar  $T$  en  $R_T$ ;
16:     Remover  $T$  de  $availTasks$ ;
17:     Actualizar  $EAT(R_T)$ ;
18: end while
19: end procedure

```

A.1.5. Sufragio

Input: Un grafo de flujo de trabajo $w = (\mathcal{V}, \mathcal{E})$

Output: Una calendarización f

```

1: while  $\exists v \in \mathcal{V}$  no completada do
2:      $t \leftarrow$  Obtener conjunto de tareas listas, no calendarizadas, con padres calendarizados;
3:     SCHED( $tasks$ );
4: end while
5: procedure SCHED( $availTasks$ )
6:     while  $\exists t \in availTasks$  no calendarizadas do
7:         for  $t \in availTasks$  do
8:              $res \leftarrow$  Obtener recursos disponibles para  $t$ ;
9:             for  $r \in res$  do
10:                 Calcular  $ECT(t, r)$ ;
11:             end for
12:              $R_t^1 \leftarrow \arg \min_{r \in res} ECT(t, r);$ 
13:              $R_t^2 \leftarrow \arg \min_{r \in res, r \neq R_t^1} ECT(t, r);$ 
14:              $suf_t \leftarrow ECT(t, R_t^2) - ECT(t, R_t^1);$ 
15:         end for
16:          $T \leftarrow \arg \max_{t \in availTasks} suf_t;$ 
17:         Calendarizar  $T$  en  $R_T^1$ ;
18:         Remover  $T$  de  $availTasks$ ;
19:         Actualizar  $EAT(R_T)$ ;
20:     end while
21: end procedure

```

A.1.6. HEFT

Input: Un grafo de flujo de trabajo $w = (\mathcal{V}, \mathcal{E})$

Output: Una calendarización f

- 1: Calcular *Tiempo de Ejecución Promedio* (1) para cada tarea $v \in \mathcal{V}$;
- 2: Calcular *Tiempo de Transferencia de Datos Promedio* (2) entre tareas y sus sucesores;
- 3: Calcular *Rango* para cada tarea, de acuerdo a (3) (4);

```

4: Ordenar las tareas por Rango, en orden decreciente en una lista  $Q$ ;
5: while  $Q \neq \emptyset$  do
6:    $t \leftarrow$  Remover la primera tarea de  $Q$ ;
7:    $r \leftarrow$  Encontrar un recurso que pueda ejecutar  $r$  en el menor tiempo;
8:   Calendarizar  $t$  en  $r$ ;
9: end while

```

A.1.7. Híbrido

Input: Un grafo de flujo de trabajo $w = (\mathcal{V}, \mathcal{E})$

Output: Una calendarización f

```

1: Calcular Peso de cada tarea y arista, de acuerdo a (1)(2);
2: Calcular Rango para cada tarea, de acuerdo a (3) (4);
3: Ordenar las tareas por Rango, en orden decreciente en una lista  $Q$ ;
4:  $i \leftarrow 0$ 
5: Crear un grupo  $G_i$ ;
6: while  $Q \neq \emptyset$  do
7:    $t \leftarrow$  Remover la primera tarea de  $Q$ ;
8:   if  $t$  tiene una dependencia con una tarea en  $G_i$  then
9:      $i \leftarrow i + 1$ ;
10:    Crear un grupo  $G_i$ ;
11:   end if
12:   Agregar  $t$  a  $G_i$ ;
13: end while
14:  $j \leftarrow 0$ ;
15: while  $j \leq i$  do
16:   Calendarizar tareas en  $G_i$  usando un algoritmo batch;
17:    $j \leftarrow j + 1$ ;
18: end while

```

A.1.8. TANH

Input: Un grafo de flujo de trabajo $w = (\mathcal{V}, \mathcal{E})$

Output: Una calendarización f

```

1: Calcular Parámetros para cada nodo tarea;
2: Agrupar tareas del flujo de trabajo;
3: if Número de clusters  $\geq$  Número de recursos disponibles then
4:   Reducir el Número de clusters al Número de recursos disponibles;
5: else
6:   Ejecutar duplicación de tareas;
7: end if

```

A.2. Algoritmos de Calidad en el Servicio

A.2.1. BackTracking

A.2.2. Deadline Distribution

A.2.3. CTC

A.2.4. Deadline-MDP

A.2.5. LOSS/GAIN

A.3. Metaheurísticos

A.3.1. Algoritmo Genético

A.3.2. GRASP

Input: Un grafo de flujo de trabajo $w = (\mathcal{V}, \mathcal{E})$

Output: Una calendarización f

```
1: while Criterio de terminación no satisfactorio do
2:    $sched \leftarrow \text{CREATE\_SCHEDULE}(w)$ ;
3:   if  $sched$  es mejor que  $bestSched$  then
4:      $bestSched \leftarrow sched$ ;
5:   end if
6: end while
7: procedure  $\text{CREATE\_SCHEDULE}(workflow)$ 
8:    $solution \leftarrow \text{CONSTRUCT\_SOLUTION}(workflow)$ ;
9:    $nSolution \leftarrow \text{LOCAL\_SEARCH}(solution)$ ;
10:  if  $nSolution$  es mejor que  $solution$  then
11:    return  $nSolution$ ;
12:  end if
13:  return  $solution$ ;
14: end procedure
15: procedure  $\text{CONSTRUCT\_SOLUTION}(workflow)$ 
16:  while calendarización no completada do
17:     $T \leftarrow$  Obtener tareas listas sin asignar;
18:    Crear RCL para cada  $t \in T$ ;
19:     $subSolution \leftarrow$  Seleccionar un recurso aleatoriamente para cada  $t \in T$  de su RCL;
20:     $solution \leftarrow solution \cup subSolution$ ;
21:    Actualizar información para futuros RCL;
22:  end while
23:  return  $solution$ ;
24: end procedure
```

```
25: procedure LOCALSEARCH(solution)
26:   nSolution  $\leftarrow$  Encontrar una solución local óptima;
27:   return nSolution;
28: end procedure
```

A.3.3. Recocido simulado

Bibliografía

- [1] Cluster Híbrido de Supercómputo — Xiuhcoatl — Cinvestav. <http://clusterhibrido.cinvestav.mx/>. Consultado el 10 de noviembre de 2013.
- [2] Laboratorio de Supercómputo y Visualización en Paralelo. <http://supercomputo.izt.uam.mx/infraestructura/aitzaloa.php>. Consultado el 10 de noviembre de 2013.
- [3] Sebastian Alcerreca Alcocer. Emisión de Comprobantes Fiscales Digitales (CFD) - Desarrollo del sistema de emisión, distribución y almacenamiento de CFD de la corporación Montecito. Tesis de licenciatura, Instituto Tecnológico Autónomo de México, 2013.
- [4] Rajkumar Buyya, Chee Shin Yeo, Srikumar Venugopal, James Broberg, and Ivona Brandic. Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. *Future Generation computer systems*, 25(6):599–616, 2009.
- [5] Ewa Deelman, Scott Callaghan, Edward Field, Hunter Francoeur, Robert Graves, Nitin Gupta, Vipin Gupta, Thomas H Jordan, Carl Kesselman, Philip Maechling, et al. Managing large-scale workflow execution from resource provisioning to provenance tracking: The cybershake example. In *e-Science and Grid Computing, 2006. e-Science'06. Second IEEE International Conference on*, pages 14–14. IEEE, 2006.
- [6] J Octavio Gutierrez-Garcia and Kwang Mong Sim. Agent-based cloud workflow execution. *Integrated Computer-Aided Engineering*, 19(1):39–56, 2012.
- [7] Charles E Leiserson, Ronald L Rivest, Clifford Stein, and Thomas H Cormen. *Introduction to algorithms*. The MIT press, 2001.
- [8] Muthucumaru Maheswaran, Shoukat Ali, HJ Siegal, Debra Hensgen, and Richard F Freund. Dynamic matching and scheduling of a class of independent tasks onto heterogeneous computing systems. In *Heterogeneous Computing Workshop, 1999.(HCW'99) Proceedings. Eighth*, pages 30–44. IEEE, 1999.
- [9] Michael Mair, Jun Qin, Marek Wieczorek, and Thomas Fahringer. Workflow conversion and processing in the ASKALON grid environment. In *2nd Austrian Grid Symposium*, pages 67–80. Citeseer, 2007.

- [10] Angela O'Brien, Steven Newhouse, and John Darlington. Mapping of scientific workflow within the e-protein project to distributed resources. In *UK e-Science All Hands Meeting*, pages 404–409, 2004.
- [11] Krithi Ramamritham, John A. Stankovic, and P-F Shiah. Efficient scheduling algorithms for real-time multiprocessor systems. *Parallel and Distributed Systems, IEEE Transactions on*, 1(2):184–194, 1990.
- [12] Jeffrey D. Ullman. NP-complete scheduling problems. *Journal of Computer and System sciences*, 10(3):384–393, 1975.
- [13] Wil MP van Der Aalst, Arthur HM Ter Hofstede, Bartek Kiepuszewski, and Alistair P Barros. Workflow patterns. *Distributed and parallel databases*, 14(1):5–51, 2003.
- [14] Marek Wieczorek, Andreas Hoheisel, and Radu Prodan. Taxonomies of the multi-criteria grid workflow scheduling problem. In *Grid Middleware and Services*, pages 237–264. Springer, 2008.
- [15] Jia Yu, Rajkumar Buyya, and Kotagiri Ramamohanarao. Workflow scheduling algorithms for grid computing. In *Metaheuristics for scheduling in distributed computing environments*, pages 173–214. Springer, 2008.