

Universidad Central “Marta Abreu” de Las Villas
Facultad de Matemática, Física y Computación



TRABAJO DE DIPLOMA

DISTRIBUIDOR DE FRAGMENTOS FÍSICOS DE DATOS EN SQL SERVER

Autor: Albanis Peña Pacheco
Tutores: Dr. Abel Rodríguez Morffi
Dra. Luisa M. González González

Curso 2011-2012

Hago constar que el presente trabajo fue realizado en la Universidad Central “Marta Abreu” de Las Villas como parte de la culminación de los estudios de la especialidad de Ciencias de la Computación, autorizando a que el mismo sea utilizado por la institución, para los fines que estime conveniente, tanto de forma parcial como total y que además no podrá ser presentado en eventos ni publicado sin la autorización de la Universidad.

Firma del autor

Los abajo firmantes, certificamos que el presente trabajo ha sido realizado según acuerdos de la dirección de nuestro centro y el mismo cumple con los requisitos que debe tener un trabajo de esta envergadura referido a la temática señalada.

Firma del tutor

Firma del Jefe del Seminario

Dedicatoria

Este trabajo lo quiero dedicar a todas las personas que de un modo u otro han jugado un rol importante en mi vida tanto dentro como fuera de la docencia, en especial a mi familia: a mi madre, a mis abuelos Santiago e Inés, a mis tíos Michel y Yaneisys, a mi familia de Placetas y a mis amigos que más que amigos han sido como hermanos para mí: Héctor, Polo, Chinong, David, Dustin, Adam, Sandro, Abdel, Raúl, Marelis, Donna, Raiza, Luis, Royne, Igor, Sergio, Julio, Maikel y a otros que se me puedan olvidar, a todos ustedes un millón de **GRACIAS**.

Agradecimientos

A toda mi **familia** por haber sido el fuerte soporte que nunca me dejó caer.

A mi hermano **Héctor** por haber estado ahí siempre que lo necesité.

A mi amigo **Polo** por su inestimable ayuda.

A mi tutor **Abel** por haberme guiado en esta recta final.

A **Enrique** por haber sido mi consultor cuando tuve dudas.

A todos mis **amigos**.

Pensamiento

“Si buscas resultados distintos, no hagas siempre lo mismo.”

Albert Einstein

Resumen

El procesamiento de bases de datos distribuidas consiste en trabajar con bases de datos, en las cuales la ejecución de transacciones, la recuperación y actualización de los datos, acontece usualmente en dos o más computadoras independientes, por lo general separadas geográficamente aunque conceptualmente basta una computadora y un acceso costoso a los datos para utilizar bases de datos distribuidas. Actualmente existen varios problemas relacionados con el diseño de bases de datos distribuidas ya que se han propuesto resultados teóricos que son difíciles de utilizar o no han solucionado totalmente el problema para el que se han desarrollado. Uno de ellos es la obtención de esquemas físicos de bases de datos distribuidas. Esto se debe fundamentalmente a la fuerte dependencia que tiene de los sistemas gestores de bases de datos.

Este trabajo muestra la creación de una nueva versión de la herramienta DISTRIBUTOR que considera un detallado examen de las características que debe seguir un esquema de replicación para Microsoft SQL Server y ofrece una solución computacional al problema de generación de esquemas físicos mediante replicación a través de la creación de scripts en Transact-SQL, tomando como entrada un documento XML que contiene el resultado del diseño de una base de datos distribuida obtenido por la herramienta SIADBDD. Esta aplicación fue desarrollada en Visual C # como parte de Visual Studio 2010.

Abstract

Distributed database is a collection of data that are connected by logic, stored separately and used together in a computer network. In other words, the distribution of this data allows other units can access data from a particular unit. In Distributed Database there are advantages and disadvantages.

There are several problems concerning to the design a distributed database. Some theoretical outcomes have been proposed but most of them are difficult to implement or to use, and none of them fully solve the problem for which they have created. One of these problems is the final allocation of data in which each fragment of data must be allocated to one or several sites, where the fragment will be saved. This is mainly due to the strong dependence of the database management systems.

The present work presents a new version of the tool DISTRIBUTOR which considers a detailed examination of the issues must follow a replicated schema generation in Microsoft SQL Server and provides a computational solution to the problem of physical schema generation by replication through the creation of scripts in Transact-SQL taking as input a XML document containing the SIADBDD tool design results obtained. This application was developed in Visual C # as part of Visual Studio 2010.

Tabla de Contenidos

Introducción	1
Capítulo 1. Diseño de Bases de Datos Distribuidas	5
1.1 Generalidades	5
1.2 Diseño lógico.....	7
1.2.1 Fragmentación.....	7
1.2.2 Ubicación	9
1.3 Diseño físico de distribución de datos.....	10
1.4 Generación de esquemas distribuidos en SQL Server	10
1.4.1 Distribución y replicación	11
1.4.2 Replicación en Microsoft SQL Server.....	11
1.5 Consideraciones parciales.....	18
Capítulo 2. Concepción de una herramienta para la distribución física de datos en Microsoft SQL Server.....	19
2.1 Diseño de la herramienta DISTRIBUTOR.....	19
2.1.1 Diagrama de Clases.....	19
2.1.2 Diagrama de Secuencia	21
2.2 Implementación	23
2.3 Generación de scripts en Microsoft SQL Server 2008	26
2.4 Conclusiones parciales.....	30
Capítulo 3. Herramienta para generar esquemas físicos distribuidos en Microsoft SQL Server ..	31
3.1 Presentación de DISTRIBUTOR	31

3.2	Descripción de un caso de estudio	37
3.3	Solución.....	39
3.4	Consideraciones parciales.....	42
	Conclusiones.....	43
	Recomendaciones.....	44
	Referencias Bibliográficas.....	45
	Anexos	48

Introducción

Las innovaciones tecnológicas recientes y el abaratamiento del hardware han estimulado el desarrollo de los sistemas de bases de datos distribuidas (SBDD). Según Ózsu y Valduriez (Ózsu and Valduriez, 1999), una base de datos distribuida (BDD) es una colección de múltiples bases de datos (BD), lógicamente interrelacionadas distribuidas sobre una red de computadoras. Así, el objetivo fundamental de los SBDD es integrar la manipulación de datos para que sean presentados al usuario como una única colección de datos global y coherente.

La distribución de datos involucra el hecho de que los datos no residen necesariamente en el mismo sitio, pero poseen propiedades comunes que los vinculan, y se facilita su acceso a través de una interfaz común (Ceri et al., 1982). Es necesario destacar que los enlaces entre los datos se llevan a cabo en una red de comunicación, lo que implica generalmente que los sitios estén localizados en diferentes áreas geográficas y con capacidad de procesamiento autónomo.

Técnicamente, los sistemas distribuidos se adaptan mejor a las necesidades de las organizaciones descentralizadas, porque reflejan más adecuadamente su estructura y tienen importantes ventajas con respecto a los centralizados. La descentralización se justifica desde el punto de vista tecnológico ya que permite autonomía local y promueve la evolución de los sistemas, así como los cambios en los requerimientos de los usuarios, proporciona una arquitectura de sistemas simple, flexible y tolerante a fallos, y ofrece buenos rendimientos. Estas razones a favor de la descentralización también añaden nuevas complejidades a su diseño e implementación.

Las diferentes técnicas de distribución de datos generalmente se basan en la semántica de ellos, y se rigen por principios idénticos que las BD centralizadas, incorporando otros detalles particulares, como la fragmentación de las entidades y su posterior ubicación en los diferentes sitios de la red. La fragmentación es muy útil para mejorar los tiempos de respuesta y garantizar el paralelismo de un sistema (Baião et al., 2004, Coulon et al., 2005, Johansson et al., 2000, Lin et al., 2005).

En el diseño de la distribución influyen muchos factores relacionados con la BD, las aplicaciones que acceden a la misma, la comunicación de la red, y sobre el sistema de computadoras que la soporta; lo que hace que sea muy complicada la formulación de un problema de distribución. Por tanto es necesario decidir cómo fragmentar y distribuir los datos sobre los diferentes sitios y cuáles de estos datos deben ser replicados. Además, el diseño físico de una BDD exige decisiones y procesamiento complejos respecto a la ubicación de los fragmentos de datos.

Como **antecedentes** a este trabajo se tiene que en el Grupo de Bases de Datos del Centro de Estudios de Informática de la Universidad Central “Marta Abreu” de Las Villas se desarrollaron herramientas de ayuda al diseño de BDD integradas en SIADBDD que es un ambiente único de coordinación entre ellas a través de un catálogo compartido. DISTRIBUTOR es una de estas herramientas la cual se encarga de la asignación física de fragmentos a los diferentes sitios de procesamiento identificados en la red donde se ubicará la BDD que se diseñe. Esta herramienta realiza la ubicación mediante la creación y posterior ejecución de un conjunto de scripts en cada uno de los sitios. Estos scripts usan en repertorio de instrucciones de Transact-SQL, dialecto de SQL propio del sistema de gestión de bases de datos (SGBD) Microsoft SQL Server, de gran popularidad en Cuba, pero ignora muchos de los elementos relacionados con la configuración del ambiente de replicación, aspecto necesario para una correcta orquestación de este ambiente. Consecuentemente el **problema de la investigación** que se presenta para este trabajo está dado por la necesidad de obtener una solución al problema de ubicación física permitiendo la configuración completa del ambiente de replicación soportado por SQL Server que se integre dinámicamente a SIADBDD.

Como **hipótesis** se tiene que una herramienta para la generación de esquemas físicos de BDD mediante el ambiente de replicación soportado por Microsoft SQL Server que no dependa directamente de las estructuras internas de SIADBDD, facilita el trabajo de los diseñadores de BDD y simplifica el mantenimiento posterior.

El **objetivo** de este trabajo es crear una herramienta para la generación de los fragmentos físicos de una BDD que incluya la configuración completa del ambiente de

replicación soportado en Microsoft SQL Server, con el fin de obtener una solución independiente que se integre dinámicamente con SIADBDD.

Como **objetivos específicos** se tiene los siguientes:

- Identificar los aspectos de programación en Transact-SQL que permiten crear entornos de replicación de datos para materializar diseños de BDD a nivel físico en Microsoft SQL Server.
- Diseñar una herramienta para generar esquemas físicos de ubicación usando los elementos programáticos de replicación identificados, que use como entrada un documento XML con la información obtenida por SIADBDD y los datos adicionales requeridos para la configuración del entorno de replicación en Microsoft SQL Server.
- Crear una herramienta en forma de asistente según el diseño obtenido que resuelva los problemas de la versión anterior de DISTRIBUTOR versión anterior respecto a la posibilidad de configurar la replicación multimaster, personalizar la creación de las bases de datos locales y reconocer la duplicación de tablas aun cuando no se apliquen filtros.
- Probar la herramienta con un caso de estudio real.

Las **preguntas de la investigación** son las siguientes:

- ¿Qué elementos obtenidos como resultado del diseño de BDD en SIADBDD son necesarios para realizar la ubicación física de los fragmentos?
- ¿Qué características debe poseer un asistente para ubicación física de fragmentos de datos en Microsoft SQL Server que resulte favorable al diseñador?
- ¿Qué aspectos de la replicación de datos se deben considerar al ubicar fragmentos físicos en SQL Server?

Relacionado con los antecedentes antes planteados, se puede decir que este estudio se **justifica** por su importancia desde el punto de vista práctico en la integración que tendrán los resultados esperados en una herramienta de ayuda al diseño de BDD.

Este trabajo es **viable**, ya que en el laboratorio de Bases de Datos existen posibilidades de realización de esta investigación, porque se cuenta con recursos materiales para ello, expertos para la consulta, y oportunidades de utilización de los resultados.

Los resultados obtenidos tienen **valor práctico** porque la herramienta constituye una solución independiente que se integra dinámicamente con SIADBDD para lograr completar diseños de BDD en Microsoft SQL Server.

La tesis está **estructurada** en tres capítulos de la siguiente forma:

El capítulo 1 trata sobre aspectos generales del diseño de BDD y particulariza en la generación de esquemas físicos y en las etapas fundamentales de la replicación.

En el capítulo 2 se propone una solución concreta al problema planteado, mostrando su diseño y los aspectos más importantes de su implementación.

En el capítulo 3 se presenta un caso de estudio real y se muestran los resultados ofrecidos por la herramienta DISTRIBUTOR para dar solución a la problemática anterior.

Este documento culmina con las conclusiones, referencias bibliográficas y los anexos.

Capítulo 1. Diseño de Bases de Datos Distribuidas

En el diseño de BDD se debe considerar el problema de cómo distribuir los datos entre los diferentes sitios de procesamiento. Existen razones organizacionales las cuales determinan en gran medida lo anterior. Sin embargo, cuando se busca eficiencia en el acceso a la información, se deben abordar los problemas de cómo fragmentar los datos y cómo asignar cada fragmento entre los diferentes sitios de la red. En el diseño de BDD también es importante considerar si la información está replicada y cómo mantener la consistencia de los mismos. En este capítulo se tratan aspectos generales relacionados con el diseño de distribución de datos así la generación de esquemas y la replicación como una etapa esencial dentro de la creación de BDD que sirven de marco teórico referencial a este trabajo.

1.1 Generalidades

La utilización de las nuevas tecnologías así como el uso intensivo de Internet ha traído consecuencias importantes en las comunicaciones así como en la distribución de datos en localidades geográficamente separadas. Ante esta situación se buscan alternativas para integrar y compartir la información necesaria, por lo que surgen nuevas tecnologías capaces de gestionar de manera estable toda esta la información. En muchas organizaciones geográficamente distribuidas, las vías centralizadas no representan una opción factible, y la migración hacia BDD resulta natural. Muchos autores recalcan el principio de que un sistema distribuido muestra la estructura de la BD como un espejo de la estructura de la organización, con lo cual se incrementa la localidad de referencia y se reduce drásticamente el tráfico en la red (Date, 2000).

Muchas han sido las acepciones para lograr la definición de un sistema distribuido, pero en ocasiones no convergen entre sí. La más viable en correspondencia con los principios básicos perseguidos en esta investigación es la siguiente: Un sistema distribuido es una colección de computadoras independientes, que aparecen ante los usuarios del sistema como una única computadora. El logro de elevados grados de

independencia entre los programas de aplicación y los aspectos internos ha sido el resultado del perfeccionamiento, avance y progreso de las tecnologías de BD.

Como se expresó en la introducción a este trabajo, una BDD puede verse como una colección de datos que pertenecen lógicamente a un solo sistema, pero se encuentra físicamente esparcida en varios sitios de la red, es decir, se tiene un conjunto de sitios conectados entre sí mediante algún tipo de red de comunicaciones en el cual cada sitio es un sistema de BD en sí mismo, y trabajan juntos con el fin de que un usuario pueda obtener acceso a los datos de cualquier punto de la red, como si estuvieran almacenados en el propio sitio del usuario.

El diseño de la BDD tiene como objetivo lograr un mejor desempeño; y requiere de su propia teoría, metodología, y herramientas de apoyo. El diseño de distribución incluye fragmentación y ubicación replicada o no de los fragmentos (Ceri et al., 1982). El término fragmentación se refiere a la división del esquema de una BDD de acuerdo con algún criterio, mientras que la replicación está relacionada con la existencia de copias de los fragmentos en varios sitios.

Los dos principios básicos de las BDD (Ceri et al., 1987, Ma et al., 2005, Özsu and Valduriez, 1999) son reducir el intercambio de datos entre los sitios y eliminar datos irrelevantes en la ejecución de solicitudes. A la hora de abordar el problema de diseño de BDD existen dos enfoques básicos, el enfoque descendente y el ascendente. El primero es más apropiado para aplicaciones nuevas y para sistemas homogéneos y consiste en partir desde el análisis de requerimientos para definir el diseño conceptual y las vistas de usuario. A partir de ellas se define un esquema conceptual global y los esquemas externos necesarios. Se prosigue con el diseño de la fragmentación de la BD, y de aquí se continúa con la localización de los fragmentos en los sitios, creando los fragmentos físicos. Esta aproximación se completa ejecutando, en cada sitio, el diseño físico de los datos que se localizan en éste. Por tanto, varios autores coinciden en que el proceso de diseño de BDD debe ser dividido en cuatro pasos fundamentales (Baião et al., 2002, Bellatreche et al., 2000, Ceri et al., 1987, Hababeh et al., 2004, Navathe et al., 1995, Özsu and Valduriez, 1999, Schewe, 2002): (1) Diseño del esquema conceptual, (2) Diseño de la fragmentación, (3) Diseño de la asignación, y (4)

Diseño de la BD física. Los problemas 1 y 4 son comunes a las BD centralizadas (BDC) y a las BDD, mientras que los problemas 2 y 3 caracterizan el diseño de BDD. El segundo enfoque se utiliza específicamente cuando se parte de BD existentes para obtener BDD integradas. En forma resumida este diseño ascendente de BDD requiere de la selección de un modelo de bases de datos común para describir el esquema global de la base de datos.

La tarea de lograr todas las funcionalidades en los SBDD tiene una alta complejidad, pero encontrar soluciones óptimas es aún más complejo. Considerar el diseño de la distribución de datos, decidir cómo fragmentar y distribuir los datos sobre los diferentes sitios y cuáles de estos datos deben ser replicados son grandes retos que existen.

1.2 Diseño lógico

El modelo lógico de las BDD fue formulado inicialmente por Ceri et al. en 1983 (Ceri et al., 1983), actualmente considerado como un modelo clásico consolidado. El mismo incluye el problema de cómo fragmentar y distribuir los datos óptimamente en los sitios de una red. Este problema tiene como principio clave alcanzar máxima localidad de los datos, ubicándolos tan cerca como sea posible de las aplicaciones que los utilizan, lo que permite reducir el tráfico de comunicaciones en la red. El diseño lógico se convierte en parte en la especificación funcional que se usa en el diseño físico. El diseño lógico es independiente de la tecnología.

Este trabajo está acorde con la estrategia descendente de diseño de BDD porque es adecuada para el desarrollo inicial de un sistema de BDD sin tener restricciones de otros sistemas ya instalados y que deban ser integrados de alguna manera al sistema distribuido como corresponde a la estrategia ascendente. Bajo el enfoque descendente, en una primera fase los esquemas globales se dividen en subconjuntos llamados fragmentos; y en una segunda fase los fragmentos son ubicados en los diferentes sitios.

1.2.1 Fragmentación

Un sistema soporta fragmentación de datos si es posible dividir sus relaciones en fragmentos. Existen dos formas básicas de realizar esta división: horizontalmente o

verticalmente. Existe otro tipo no básico de fragmentación que es generada por la combinación de fragmentos horizontales y verticales, denominada fragmentación mixta. Aun cuando esta fragmentación no se considere un tipo primitivo, en muchas situaciones reales resulta imprescindible disponer de ella.

La fragmentación mixta puede realizarse de tres formas diferentes:

1. Desarrollando primero la fragmentación vertical y luego aplicar la fragmentación horizontal sobre los fragmentos verticales (llamada partición VH).
2. Desarrollando primero una fragmentación horizontal y luego aplicar la fragmentación vertical sobre los fragmentos horizontales (llamada partición HV).
3. Aplicando sobre una relación, de forma simultánea y no secuencial la fragmentación horizontal y vertical, generando una rejilla y los fragmentos formaran las celdas de dicha rejilla.

Una decisión compleja e importante que afecta el rendimiento de las aplicaciones es determinar cuál relación debe ser fragmentada. En este sentido influyen determinados parámetros que deben caracterizar tanto a las aplicaciones como a la BD, pero es necesario resaltar que el grado de la fragmentación puede ir de un extremo en que no se fragmenta, a otro donde se fragmenta hasta obtener atributos individuales o tuplas, para el caso de la fragmentación vertical y horizontal respectivamente. Los efectos de fragmentar unidades muy grandes pueden influir negativamente en relación con la réplica. Si el fragmento no es replicado se incrementan desmesuradamente los accesos remotos, y si es replicado en todos los sitios causa serios problemas con la actualización de los datos y es necesario incrementar la capacidad de almacenamiento de los sitios de procesamiento. Por otra parte, si es necesario recuperar datos ubicados en más de un fragmento, la fragmentación de unidades muy pequeñas provoca el uso de acoples muy costosos, y además, el problema de la localización puede resultar extremadamente complejo. Por tal motivo, el diseñador debe encontrar un nivel de fragmentación adecuado, que será un compromiso entre ambos extremos.

El propósito del diseño de la fragmentación es determinar los fragmentos que constituyen unidades lógicas de asignación y consiste en agrupar tuplas o atributos con las mismas propiedades. Los fragmentos deben ser colecciones homogéneas de

información desde el punto de vista de accesos de las aplicaciones, o sea, que todas las instancias de los fragmentos sean accedidas uniformemente por las aplicaciones que las usan.

El diseño de la fragmentación antecede al proceso de asignación de fragmentos a sitios de procesamiento, reflejando el criterio de mantener los datos locales en el sitio donde frecuentemente son accedidos por las aplicaciones; es por ello que los fragmentos constituyen una unidad apropiada de asignación. Si las aplicaciones que manipulan un esquema están ubicadas en sitios diferentes, se pueden seguir dos alternativas con todo el esquema como unidad de distribución: replicado o no replicado (Bhalla and Hasegawa, 2005, Gançarski et al., 2002, Johansson et al., 2000).

Muchos de los problemas relacionados con la fragmentación y asignación en el diseño de BDD tienen una modelación matemática compleja y son considerados de la clase NP-Completo (Baião et al., 2003, Özsu and Valduriez, 1999, Lee and Baik, 2004, Pérez et al., 2005), donde no hay garantía de encontrar una solución óptima con algoritmos determinísticos en un tiempo polinomial. Así, su complejidad es tal, que cualquier algoritmo que resuelve óptimamente cada uno de sus casos requiere un esfuerzo computacional que crece exponencialmente en función del tamaño del problema, en dependencia de la cantidad de fragmentos y de sitios.

Intentando dar solución a estos problemas se han presentado propuestas (Baião et al., 2004, Pérez et al., 2005, Savonnet et al., 1999, Tamhankar and Ram, 1998, Navathe et al., 1995, Mei and Sheng, 1992, Ceri et al., 1987, Ceri and Pernici, 1985) que sugieren el uso de heurísticas para algunos de ellos.

1.2.2 Ubicación

La ubicación de recursos a lo largo de una red de computadoras es un problema que ha sido estudiado extensivamente. Solo una pequeña proporción de los estudios realizados abarcan la distribución de datos de una BD (Bertone, 2004). El problema de la ubicación de los datos consiste en encontrar el esquema de distribución óptimo. La optimización puede ser definida en función de costo mínimo o desempeño máximo. La construcción de un modelo que permita evaluar todos los aspectos que involucra la

optimización del esquema de distribución es muy compleja. Las variables a tener en cuenta son muchas y la incidencia de cada una de ellas también lo es, como por ejemplo el costo de replicar los fragmentos a lo largo de la red (Özsu and Valduriez, 1991).

1.3 Diseño físico de distribución de datos

El diseño físico es el proceso de escoger las estructuras de almacenamiento en disco y métodos de acceso a los datos más adecuada para lograr un buen rendimiento de la BD. Para realizar el diseño físico es importante conocer la carga de trabajo, combinación de consultas y actualizaciones que debe soportar la BD y los requerimientos del usuario. También es importante que el diseñador conozca las técnicas de procesamiento de consultas e indexación soportadas por el SGBD.

El diseño físico traduce el diseño lógico en una solución implementable y económicamente efectiva. El diseño físico es el proceso de producir la descripción de la implementación de la BD en memoria secundaria: estructuras de almacenamiento y métodos de acceso que garanticen un acceso eficiente a los datos. Entre el diseño físico y el diseño lógico hay una retroalimentación, ya que alguna de las decisiones que se tomen durante el diseño físico para mejorar las prestaciones pueden afectar a la estructura del esquema lógico.

El diseño físico está íntimamente ligado a una alternativa tecnológica. Ante la acelerada evolución tecnológica es importante considerar los estándares del momento y las tendencias ya que una mala decisión tendrá, inevitablemente, una implicación en los costos.

1.4 Generación de esquemas distribuidos en SQL Server

A continuación se analizan los aspectos relacionados con la distribución y replicación en Microsoft SQL Server, que es uno de los SGBD más usados en Cuba y para el cual se plantea una solución de generación de esquemas físicos mediante replicación.

1.4.1 Distribución y replicación

Como se ha podido entender, los datos se distribuyen para lograr mejoras del desempeño puesto que cada servidor solamente estará gestionando los datos asociados con ese servidor y las relaciones son tan compactas como sea posible o por reducción de costos donde muchos servidores pequeños son menos costosos de mantener que un único servidor monolítico. Un único servidor representa un único punto de fracaso. Si un servidor falla en un escenario de servidores múltiples, los otros servidores pueden continuar con el servicio en sus ubicaciones.

Por otra parte, los datos se replican para apoyar la distribución de otros datos ya que en muchas ocasiones los datos distribuidos pueden contener llaves de relaciones replicadas, para mantener la integridad referencial donde las relaciones de apoyo deben ser replicadas al sitio de datos distribuidos, para permitir la accesibilidad a los mismos datos porque puede ser deseable replicar resúmenes de información a cada servidor y estas estadísticas pueden ser generadas desde cualquier servidor.

1.4.2 Replicación en Microsoft SQL Server

La replicación de datos permite que cierta información de la BD sea almacenada en más de un sitio y consiste en el transporte de esta entre dos o más servidores, permitiendo que ciertos datos de la BD estén almacenados en más de un sitio, y así aumentar la disponibilidad de los mismos y mejorar el rendimiento de las consultas globales (Morell, 2004).

La replicación en Microsoft SQL Server consiste en el transporte de datos entre dos o más instancias de servidores SQL Server. Para ello, Microsoft SQL Server brinda un conjunto de soluciones que permite copiar, distribuir y posiblemente modificar datos de toda la organización. Se incluyen, además, varios métodos y opciones para el diseño, implementación, supervisión y administración de la replicación, que le ofrecen la funcionalidad y flexibilidad necesarias para distribuir datos y mantener su coherencia (Morell, 2004).

El modelo de replicación en Microsoft SQL Server toma como referencia la metáfora de la industria de las publicaciones; está formado por: publicador, distribuidor, suscriptor,

publicación, artículo y suscripción; y varios agentes responsabilizados de copiar los datos entre el publicador y el suscriptor. Estos agentes son: agente de instantáneas, agente de distribución, agente del lector del registro, agente del lector de cola y agente de mezcla (Morell, 2004). A los tipos básicos de replicación de instantáneas, transaccional y de mezcla se le incorporan opciones para ajustarse aún más a los requerimientos del usuario. A continuación se explican en detalle cada uno de estos componentes.

1.4.2.1 Componentes del modelo de replicación

En este epígrafe se analizan cada uno de los componentes del modelo de replicación. Así, el publicador es un servidor que pone los datos a disposición de otros servidores para poder replicarlos; el distribuidor es un servidor que almacena la BD de distribución y almacena los datos históricos, transacciones y metadatos; los suscriptores reciben los datos replicados. Una publicación es un conjunto de artículos de una publicación de una BD. Esta agrupación de varios artículos facilita especificar un conjunto de datos relacionados lógicamente y los objetos de BD que desea replicar conjuntamente. Un artículo de una publicación puede ser una tabla de datos, la cual puede contar con todas las filas o algunas mediante el filtrado horizontal y simultáneamente contar de todas las columnas o algunas a través del filtrado vertical, un procedimiento almacenado, una definición de vista, la ejecución de un procedimiento almacenado, una vista, una vista indizada o una función definida por el usuario. Una suscripción es una petición de copia de datos o de objetos de BD a replicar. Una suscripción define qué publicación se recibirá, dónde y cuándo. Las suscripciones pueden ser de inserción (push) o de extracción (pull); y una publicación puede admitir una combinación de suscripciones de inserción y extracción. Con una suscripción de inserción el publicador propaga los cambios al suscriptor sin petición del suscriptor, los cambios pueden ser insertados bajo demanda, continuamente o bajo un horario programado, el agente de distribución o el agente de mezcla se ejecutan en el distribuidor; se utiliza cuando las publicaciones requieren movimientos de datos en tiempo real, o cuando estos deban ser sincronizados de manera continua o bajo un repetido esquema programado, más comúnmente usada con replicación de instantáneas y replicación de transacciones. Con

la subscripciones de extracción el suscriptor solicita los cambios hechos en el publicador, este permite a los usuarios en el suscriptor determinar cuándo se sincronizarán los datos modificados, los agentes de distribución o de mezcla se ejecutan en el suscriptor. Se usa cuando los datos son sincronizados bajo demanda o a un horario programado de manera no tan continua, cuando las publicaciones tienen un gran número de suscriptores, habitualmente usado con replicaciones de mezcla. El publicador, en las subscripciones de inserción, o el suscriptor, en las subscripciones de extracción, solicitan la sincronización o distribución de datos de una suscripción.

El publicador puede disponer de una o más publicaciones, de las cuales los suscriptores se suscriben a las publicaciones que necesitan, nunca a artículos individuales de una publicación. El publicador, además, detecta qué datos han cambiado durante la replicación transaccional y mantiene información acerca de todas las publicaciones del sitio. La función del distribuidor varía según la metodología de replicación implementada. En ocasiones se configura como distribuidor el mismo publicador y se le denomina distribuidor local. En el resto de los casos el distribuidor será remoto, pudiendo coincidir en algún caso con un suscriptor. Los suscriptores además de obtener sus subscripciones, en dependencia del tipo y opciones de replicación elegidas, pueden devolver datos modificados al publicador. Además puede tener sus propias publicaciones (Morell, 2004).

1.4.2.2 Escenarios típicos de la replicación

En una solución de replicación pudiera ser necesario utilizar varias publicaciones en una combinación de metodologías y opciones. En la replicación los datos o transacciones fluyen del publicador al suscriptor pasando por el distribuidor. Por lo tanto, en su configuración mínima, una topología de replicación se compone de al menos dos o tres servidores Microsoft SQL Server, que desempeñan los tres roles mencionados.

Se pudiera contar con las siguientes variantes si se modifica la ubicación del servidor distribuidor:

1. El rol de distribuidor desempeñado por el publicador (véase la figura 1.1).

2. El rol de distribuidor desempeñado por el suscriptor (véase la figura 1.2).
3. Un servidor de distribución, independiente del publicador y del suscriptor (véase la figura 1.3).



Figura 1.1.Publicador-Distribuidor

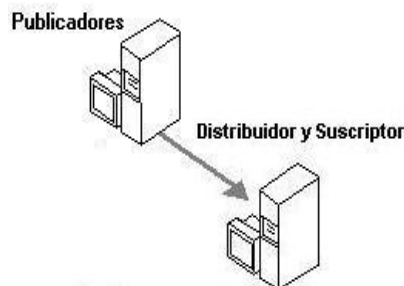


Figura 1.2. Distribuidor- Suscriptor

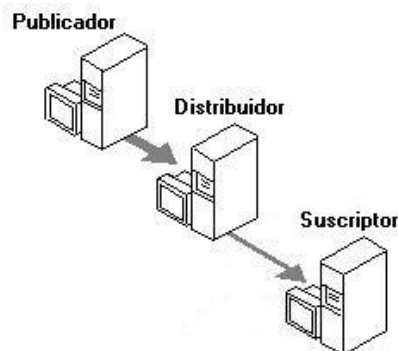


Figura 1.3.Distribuidor independiente

En la mayoría de las configuraciones, el peso fundamental de la replicación recae sobre el servidor de distribución. Por tanto, esto puede ser un criterio para determinar su ubicación, teniendo en cuenta las configuraciones o posibilidades físicas de los servidores, así como otras responsabilidades que pueden estar desempeñando: servidor de dominio, servidor de páginas Web entre otras (Morell, 2004). Existe la

posibilidad de contar con un servidor que se suscriba a una publicación y a la vez la publique para el resto de los suscriptores; esto puede ser muy útil cuando se cuente con una conexión muy costosa con el publicador principal. Evidentemente, en una configuración tal, pueden nuevamente combinarse la ubicación de los dos distribuidores y aumentar el número de variantes que pueden presentarse, pero las consideraciones para determinar la ubicación del servidor que fungirá como distribuidor son las ya mencionadas.

1.4.2.3 Tipos de replicación

Los tres tipos básicos de replicación soportados por Microsoft SQL Server son: instantáneas, transaccional y de mezcla.

1. Replicación de instantáneas: En este tipo de replicación los datos se copian tal y como aparecen exactamente en un momento determinado. Por consiguiente, no requiere un control continuo de los cambios. Las publicaciones de instantáneas se suelen replicar con menos frecuencia que otros tipos de publicaciones y puede llevar más tiempo propagar las modificaciones de datos a los suscriptores. Se recomienda utilizar este tipo de replicación cuando la mayoría de los datos no cambian con frecuencia, cuando se replican pequeñas cantidades de datos, cuando los sitios con frecuencia están desconectados y es aceptable un período de latencia largo, o sea, la cantidad de tiempo que transcurre entre la actualización de los datos en un sitio y en otro. En ocasiones se hace necesario utilizarla cuando están involucrados algunos tipos de datos (text, ntext, e image) cuyas modificaciones no se registran en el registro de transacciones y por tanto no se pueden replicar utilizando la metodología de replicación transaccional (Morell, 2004). Con la opción de actualización inmediata en el suscriptor, se permite a los suscriptores actualizar datos solamente si el publicador los va a aceptar inmediatamente. Si el publicador los acepta, se propagan a otros suscriptores. El suscriptor debe estar conectado de forma estable y continua al publicador, para poder realizar cambios en el suscriptor. Esta opción es útil en escenarios en los que tienen lugar unas cuantas modificaciones ocasionales en los servidores suscriptores.

2. Replicación transaccional: En este caso se propaga una instantánea inicial de datos a los suscriptores, y después, cuando se efectúan las modificaciones en el publicador, las transacciones individuales se propagan a los suscriptores. Microsoft SQL Server almacena las transacciones que afectan a los objetos replicados y propaga esos cambios a los suscriptores de forma continua o a intervalos programados. Al finalizar la propagación de los cambios, todos los suscriptores tendrán los mismos valores que el publicador. Este tipo de replicación suele utilizarse cuando se desea que las modificaciones de datos se propaguen a los suscriptores normalmente pocos segundos después de producirse; se necesita que las transacciones sean atómicas, que se apliquen todas o ninguna al suscriptor; los suscriptores se conectan en su mayoría al publicador; su aplicación no puede permitir un período de latencia largo para los suscriptores que reciban cambios. Esta es útil en escenarios en los que los suscriptores pueden tratar a sus datos como de sólo lectura, pero necesitan cambios a los datos con una cantidad mínima de latencia. Con el uso de la opción de actualización inmediata en el suscriptor se pierde aún más la autonomía de sitio, pero se reduce el tiempo en el cual los sitios actualizan sus copias de los datos. Para hacer modificaciones en la BD del suscriptor éstas se realizan o intentan también en la base de datos publicador en una confirmación de dos fases (2PC) por lo que si su modificación se confirma, indica que es válida y luego, en cuestión de minutos, o según la planificación hecha, estos cambios son duplicados a las demás BD suscriptoras.
3. Replicación de mezcla: Esta permite que varios sitios funcionen en línea o desconectados de manera autónoma, y mezclar más adelante las modificaciones de datos realizadas en un resultado único y uniforme. La instantánea inicial se aplica a los suscriptores. A continuación Microsoft SQL Server hace un seguimiento de los cambios realizados en los datos publicados en el publicador y en los suscriptores. Los datos se sincronizan entre los servidores a una hora programada o a petición. Las actualizaciones se realizan de manera independiente, sin protocolo de confirmación en más de un servidor; así el publicador o más de un suscriptor pueden haber actualizado los mismos datos. Por lo tanto, pueden producirse conflictos al

mezclar las modificaciones de datos. Cuando se produce un conflicto, el agente de mezcla invoca una resolución para determinar qué datos se aceptarán y se propagarán a otros sitios. Este tipo de replicación es útil cuando varios suscriptores necesitan actualizar datos en diferentes ocasiones y propagar los cambios al publicador y a otros suscriptores, cuando los suscriptores necesitan recibir datos, cuando se necesita realizar cambios sin conexión y sincronizar más adelante los cambios con el publicador y otros suscriptores y cuando la autonomía del sitio es un factor crucial. La misma es útil en ambientes en los que cada sitio hace cambios solamente en sus datos pero que necesitan tener la información de los otros sitios. Para ajustarse aún más a los requerimientos de los usuarios se incorporan opciones como son la actualización inmediata en el suscriptor, la actualización en cola y la transformación de datos replicados (Morell, 2004).

1.4.2.4 Factores para elegir el método de replicación a utilizar

En la elección de un método adecuado para la distribución de los datos en una organización influyen varios factores, los cuales se pueden agrupar en dos: factores relacionados con los requerimientos de la aplicación y factores relacionados con el entorno de red.

Dentro de los factores relacionados con los requerimientos de la aplicación, los fundamentales son (Morell, 2004): autonomía, consistencia transaccional y latencia. La autonomía de un sitio da la medida de cuanto puede operar el sitio desconectado de la base de datos publicadora. La consistencia transaccional de un sitio viene dado por la necesidad de ejecutar o no inmediatamente todas las transacciones que se han ejecutado en el servidor, o si es suficiente con respetar el orden de las mismas. La latencia de un sitio se refiere al momento en que se deben sincronizar las copias de los datos (Morell, 2004).

Entre los factores relacionados con el entorno de red están la velocidad de transmisión de datos de la red y la confiabilidad de la misma. Por otra parte en el caso que los servidores SQL Server no permanezcan todo el día encendido, como pudiera suceder

en algunas organizaciones, deben considerarse los horarios de disponibilidad de cada servidor (Morell, 2004).

La consideración de estos factores sirve de guía en la configuración del ambiente de replicación. Además deben considerarse las siguientes preguntas: ¿Qué datos se van a publicar? ¿Reciben todos los suscriptores todos los datos o sólo subconjuntos de ellos? ¿Se deben particionar los datos por sitio? ¿Se debe permitir que los suscriptores envíen actualizaciones de los datos? Y en caso de permitirlos ¿Cómo deben implementarse? ¿Quiénes pueden tener acceso a los datos? ¿Se encuentran estos usuarios en línea? ¿Se encuentran conectados mediante enlaces caros?

1.4.2.5 Fases generales para implementar y supervisar la replicación

A pesar de que existen varias formas de implementar y supervisar la replicación, el proceso de replicación es diferente según el tipo y las opciones elegidas. En general, la replicación se compone de las siguientes fases:

1. Configuración de la replicación.
2. Generación y aplicación de la instantánea inicial.
3. Modificación de los datos replicados.
4. Sincronización y propagación de los datos.

1.5 Consideraciones parciales

La configuración de la replicación en el diseño físico de BDD es un elemento importante a tener en cuenta. La topología más apropiada será configurar un servidor distribuidor-publicador y otro suscriptor. La variante replicación de mezcla es la más conveniente por las ventajas que ofrece, sobre todo porque estimula la autonomía como medida de cuánto puede operar el sitio desconectado de la BD publicadora. En este tipo de replicación se puede implementar la fragmentación a través del filtrado de datos horizontal y vertical. Se aplica el tipo de subscripción de extracción debido a que los datos serán típicamente sincronizados por petición o por una hora programada. Los suscriptores determinarán cuando se conectarán y sincronizarán los cambios.

Capítulo 2. Concepción de una herramienta para la distribución física de datos en Microsoft SQL Server

En este capítulo se abordan los aspectos relacionados con el desarrollo de la herramienta creada para dar solución al problema planteado en esta tesis. Esta se denomina DISTRIBUTOR y es la encargada de la generación de scripts para la creación de BDD mediante replicación en el SGBD Microsoft SQL Server. Para esto se utiliza la notación del lenguaje UML (acrónimo del inglés Unified Modeling Language)

2.1 Diseño de la herramienta DISTRIBUTOR

A continuación se detallan los elementos de análisis y diseño de la herramienta DISTRIBUTOR para la generación de esquemas físicos en SQL Server usando la notación de UML.

2.1.1 Diagrama de Clases

Un diagrama de clases es un esquema, patrón o plantilla para describir muchas instancias de datos posibles, que describen clases de objetos. Un diagrama de clases dado se corresponde con un conjunto infinito de diagramas de instancia. Los diagramas de clases se utilizan generalmente para mostrar clases y sus relaciones, pero también pueden utilizarse para mostrar subsistemas e interfaces.

Este tipo de diagrama muestra la relación existente entre las seis clases de este trabajo, sus nombres, atributos y métodos (véase la figura 2.1). La clase MainForm es la clase principal. Esta clase tiene servicios que utilizan los ofrecidos por las demás clases con las que está relacionada: TablasRepetidas, WizardForm, StringFunctions, XMLReader, FileHandler.

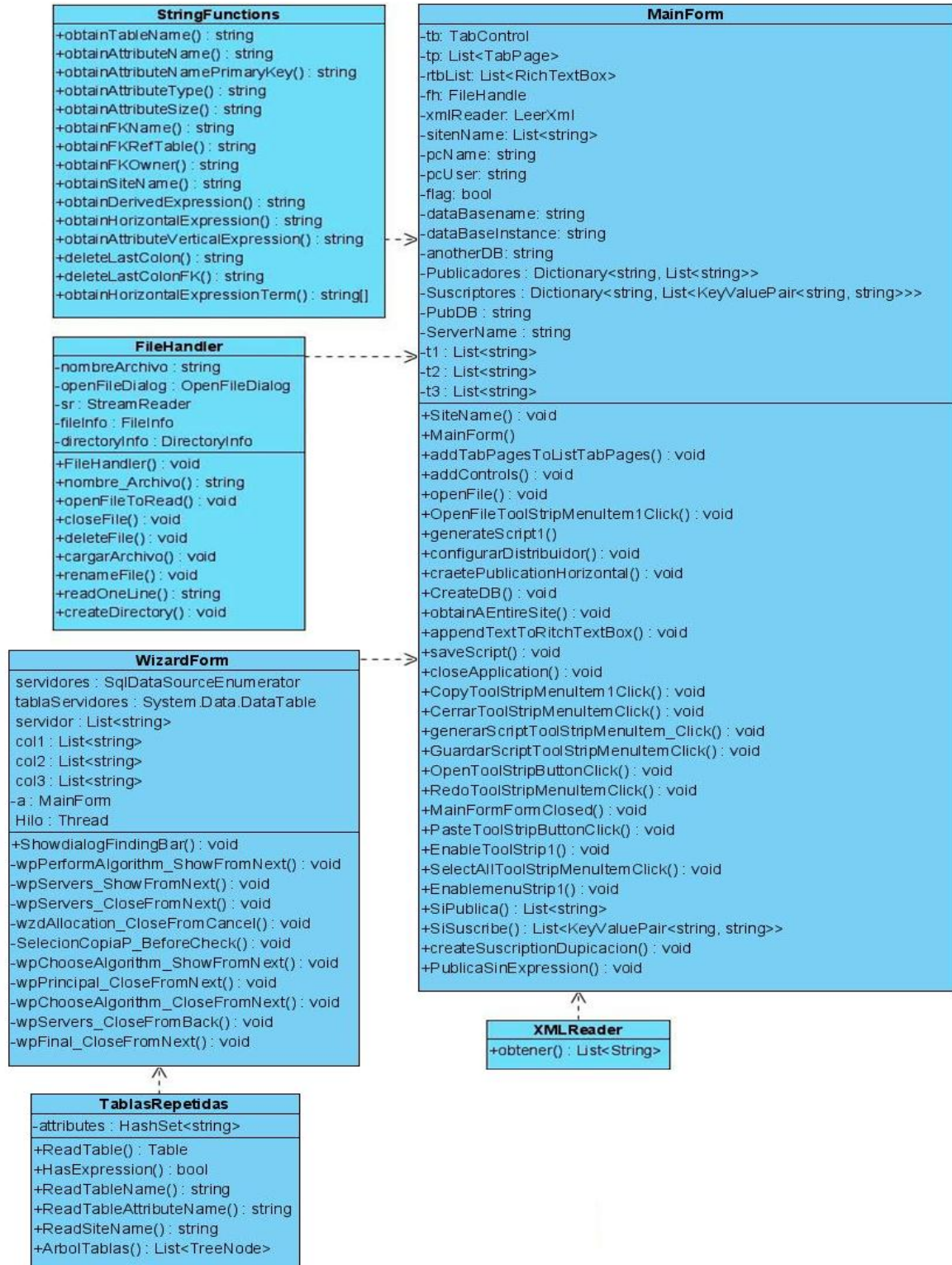


Figura 2.1. Diagrama de clases de DISTRIBUTOR.

2.1.2 Diagrama de Secuencia

El objetivo de este diagrama es indicar el orden temporal en el que se ejecutan los envíos de mensajes entre las diferentes clases. Cada línea de vida está representada por una clase. El software comienza a ejecutarse por la activación de un método de la clase MainForm que muestra la interfaz gráfica con la que el usuario va a trabajar. Primeramente se selecciona el archivo XML con el que se va a trabajar ya que esta herramienta toma como entrada una solución genérica en ese formato mediante el método openFileDialog. Este archivo es generado por servicios soportados en una biblioteca de enlace dinámico (DLL) que forma parte de SIADBDD. Esta DLL posee servicios para la creación y validación de los archivos XML que sirven de entrada a la herramienta DISTRIBUTOR que es el resultado final de esta tesis. Una vez seleccionado el XML se realiza una copia del mismo en el directorio de trabajo con nombre Catálogo.txt para el trabajo temporal. Seguidamente se muestra la forma principal MainForm con páginas para cada sitio registrado en el catálogo. Posteriormente el usuario debe elegir la opción de generar script y se muestra el asistente WizardForm el cuál obtiene información de TablasRepetidas; cuando esto se realiza se abre el archivo Catalogo.txt, se lee secuencialmente y se muestra la sentencia SQL correspondiente a esa línea.

Este método se auxilia de otros como configurarDistribuidor, createHorizontalPublication, createSuscryptionDupicacion, PublicaSinExpression. También se auxilia de StringFunctions para obtener la información relevante de cada línea. Una vez terminada la generación del script, el usuario puede determinar guardarlo en la dirección escogida por él mismo (véase la figura 2.2).

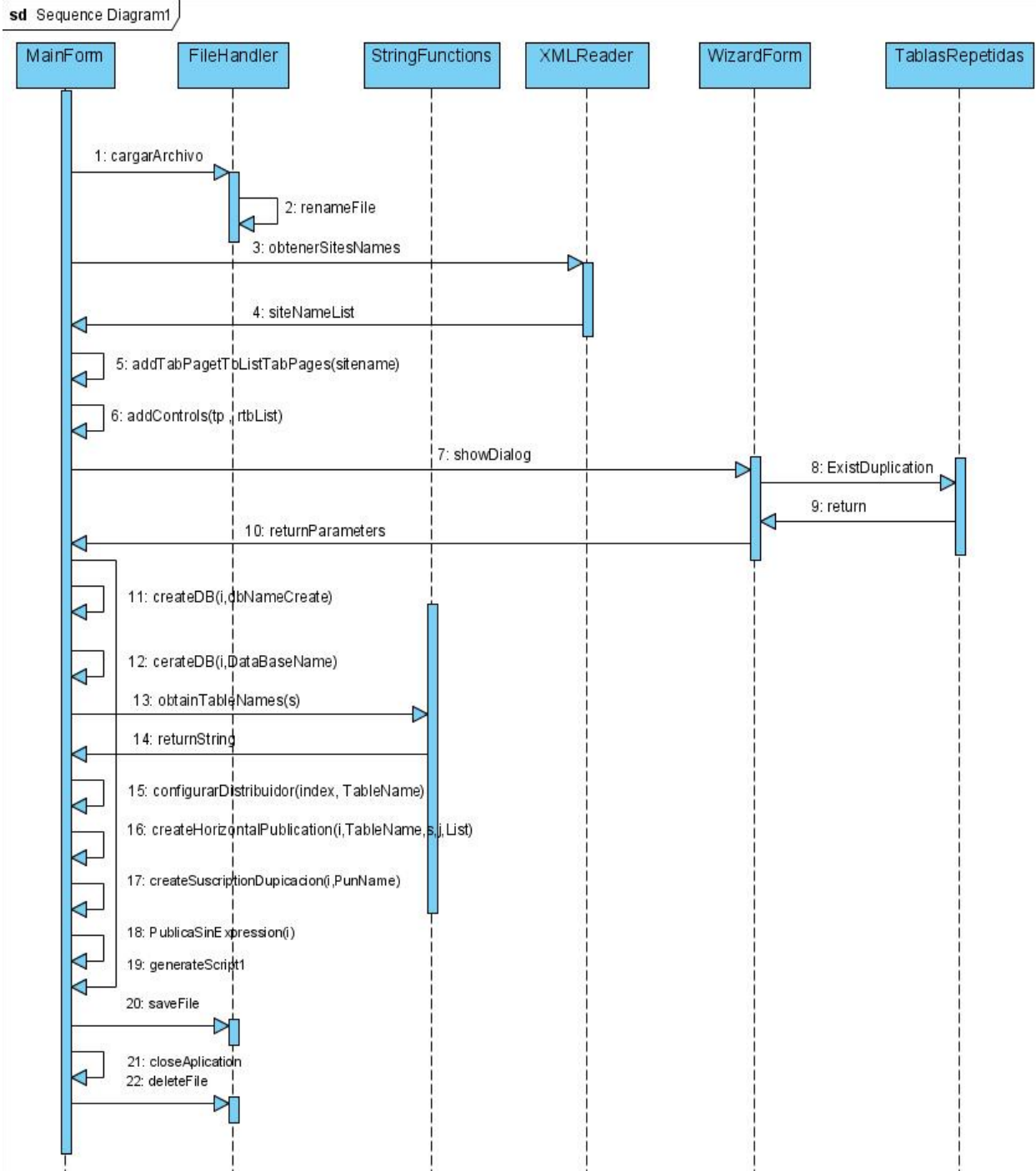


Figura 2.2. Diagrama de secuencia de DISTRIBUTOR.

2.2 Implementación

A continuación se explican brevemente las clases que están implementadas en el software, así como los métodos de cada una de ellas.

Clase WizardForm

Esta clase se utiliza para recoger los parámetros que se necesitan para la configuración de la replicación como: nombre de las instancias de los servidores, nombre de la base de datos de publicación y subscripción, también brinda la posibilidad de seleccionar la ubicación de las copias primarias para aquellas tablas que se encuentren repetidas en más de un sitio. Esta acción se realiza seleccionando el sitio que publicará los datos de modo que los demás sitios se suscribirán a éste. Esta es una característica nueva de la versión de DISTRIBUTOR que se obtiene como resultado de esta tesis. Los datos obtenidos se almacenan en la clase Datos que tiene un atributo por cada parámetro que necesita y utiliza una propiedad (property) para manejar estos valores.

Clase TablasRepetidas

Esta clase es la encargada de detectar si existe o no duplicación, para ello se auxilia de los siguientes métodos:

TableReader(): obtiene la información de una tabla.

ReadTableAttributeName(): obtiene los atributos correspondientes a cada tabla.

HasExpression(): devuelve verdadero o falso dependiendo de si la tabla en cuestión posee o no el campo <Expression>.

ArbolTablas(): devuelve un diccionario cuya llave es un String que representa el nombre de la tabla que aparece duplicada y el valor es un HashSet que contiene el nombre de los sitios que contienen dicha tabla.

Clase FileHandler

De manera general, esta clase se encarga de manipular lo referente al trabajo con archivos. Esta carga el archivo XML y crea un archivo temporal llamado catálogo.txt utilizado para la generación del script. Sus métodos son los siguientes:

openFileToRead(): crea un objeto de tipo StreamReader. El constructor de esta clase tiene un parámetro con la dirección del archivo del que se va a leer. El objeto creado es el encargado de leer el archivo.

closeFile(): cierra el archivo.

deleteFile(): elimina el archivo catálogo.txt después que deja de ser útil.

cargarArchivo(): muestra un diálogo para seleccionar el archivo XML a cargar y pone el camino seleccionado en el atributo nombreArchivo; para ello usa la propiedad nombre_Archivo.

renameFile(): hace una copia del archivo.xml pero con la extensión .txt y lo pone en el directorio de trabajo para poder crear el script a partir de este archivo.

readOneLine(): lee una línea del documento y pone el cursor en la línea siguiente.

createDirectory(): crea un directorio en Mis documentos con el nombre MSSQL.

Clase MainForm

Esta es la clase principal del software que muestra el script generado para cada sitio en una interfaz gráfica al usuario con una colección de pestañas (TabPage). Para esto se utilizan diferentes métodos. Primeramente se muestra la ventana con una lista de pestañas (TabPage) que tienen como nombre los nombres de los sitios (SitesNames) leídos del archivo XML. Esto se hace a través de los siguientes métodos:

addTabPageToListTabPage(List<string>s): crea una pestaña (TabPage) y un cuadro de texto (RichTextBox) por cada sitio; esto lo hace a partir del método obtener() de la clase XmlReader que devuelve una lista con todos los sitios del XML la cual se pasa como parámetro a este método.

addControls(List<TabPage>tp, List<RichTextBox>rt): agrega el contenedor de pestañas (TabControl), a la ventana principal (Form); pone las propiedades necesarias para mostrarla y le agrega una pestaña (TabPage) por sitio que a su vez contiene un cuadro de texto (RichTextBox), cada uno con un conjunto de propiedades para que se muestre en el contenedor de pestañas (TabControl).

openFile(): este método se invoca desde la opción Abrir archivo XML de la barra de herramientas o desde el menú Archivo. Básicamente hace una llamada a los métodos cargarArchivo(), obtener(), addTabPageToListTabPage(), addControls() y como resultado final se muestra el área de generación de script de la ventana.

generateScript1(): este método es el encargado de generar el script que se muestra en el cuadro de texto (RichTextBox); esto se hace recorriendo el archivo catálogo.txt y

escribiendo la sentencia en Transact-SQL correspondiente a cada línea leída. Este método utiliza fundamentalmente los métodos de la clase StringFunctions para obtener la información de cada línea leída y el método appendTextToRitchTextBox() que imprime en el cuadro de texto (RitchTextBox) una cadena pasada como parámetro. También utiliza otros métodos cómo configurarDistribuidor(), createHorizontalPublication(), createSuscriptionDupicacion(), PublicaSinExpression().

obtainEntireSite(string nameSite): obtiene todos las líneas de un sitio pasado por parámetro.

saveScript(): crea un archivo .txt por cada sitio con el script generado y lo guarda en la dirección seleccionada por el usuario a través de FolderBrowserDialog.

closeApplication(): cierra la aplicación y elimina el archivo catálogo.txt.

Clase StringFunctions

Esta clase cuenta con un conjunto de métodos estáticos por lo que no se necesita crear un objeto de la clase para poder utilizar sus métodos implementados para el trabajo con cadenas. Sus métodos son los siguientes:

obtainTableName(string s): devuelve el nombre de una tabla a partir de la cadena <SchemaName>nombre-de-tabla</SchemaName>

obtainAttributeName(string s): devuelve el nombre de un atributo a partir de la cadena <AttName>nombre-de-atributo</AttName>

obtainAttributeNamePrimaryKey(string s): devuelve el nombre de un atributo en la llave primaria de una tabla a partir de la cadena <AttributePK> nombre-de-atributo-llave-primaria </AttributePK>.

obtainAttributeType(string s): devuelve el tipo de un atributo a partir de la cadena<AttType>tipo-de-atributo</AttType>.

obtainAttributeSize(string s): devuelve el tamaño de un atributo a partir de la cadena <AttSize>tamaño-de-atributo</AttSize>.

obtainFKRefTable(string s): obtiene el nombre de la tabla a la que hace referencia la llave foránea a partir de la cadena <RefTable> tabla-propietaria</RefTable>.

obtainFKOwner(string s): obtiene el propietario de la llave foránea a partir de la cadena <AttributeOwnerFK>atributo-llave-tabla-propietaria</AttributeOwnerFK>.

obtainSiteName(string s): devuelve el nombre del sitio a partir de la cadena <SiteName>sitio</SiteName>.

obtainHorizontalExpression(string s): devuelve el valor de la expresión horizontal a partir de la cadena <ExpressionHorizontal>minterm</ExpressionHorizontal>.

obtainDerivedExpression(string s): devuelve el valor de la expresión derivada a partir de la cadena <ExpressionDerivada>tabla-propietaria </ExpressionDerivada>.

obtainAttributeVerticalExpression(string s): devuelve el valor del atributo de la expresión vertical a partir de la cadena <AttributeVertical>atributo-proyectado</AttributeVertical>.

obtainHorizontalExpressionTerm(string s): devuelve el valor del atributo del término de la expresión vertical.

deleteLastColon(string s): elimina la última coma de una línea en la creación de las llaves primarias. Esto es necesario acorde con la sintaxis de Transact-SQL.

deleteLastColonFK(string s): elimina la última coma de una línea en la creación de las llaves foráneas. Esto es necesario acorde con la sintaxis del Transact-SQL.

Clase XMLReader

En esta clase se lee el archivo XML y se obtienen todos los nombres de los sitios donde se ubicarán los datos. Este nombre se utiliza para crear una pestaña (TabPage) por sitio. Su método es:

obtener(stringns, string s, stringpath): este método iterativo recorre el archivo XML guardando los nombres de todos los sitios en una lista. Éste recibe como parámetro el sitio (Site), el nombre del sitio (SiteName) y el nombre del archivo. El nombre es la dirección completa incluyendo el camino.

2.3 Generación de scripts en Microsoft SQL Server 2008

La versión anterior e inicial de DISTRIBUTOR se encargaba de la asignación de fragmentos a los diferentes sitios de procesamiento identificados en la red donde se ubicaría la BDD diseñada pero ignora muchos de los elementos relacionados con la

configuración del ambiente de replicación, aspecto necesario para una correcta orquestación de este ambiente. En este acápite se explica la forma en que se generan los scripts que van a ser usados por Microsoft SQL Server 2008 para la configuración de la replicación. Para esto, la herramienta DISTRIBUTOR comienza con la lectura de un archivo XML anteriormente validado. Posteriormente este archivo se manipula con los métodos descritos anteriormente en el acápite de Implementación. Una vez que se ha obtenido la información necesaria de este archivo, se comienzan a generar los scripts y se colocan en pestañas (TabPages) con el título del nombre del sitio donde se ubicará cada script. Estos scripts crean las tablas cuyas definiciones se obtienen del archivo XML, incluyendo las definiciones de sus atributos llaves y descriptores. Una vez terminada la creación de las tablas se pasa a configurar la replicación la cual se explica a continuación.

Para configurar la replicación en Microsoft SQL Server se inicia con la configuración del servidor de distribución, luego el de publicación, y a continuación el suscriptor. Para este trabajo se escogió la arquitectura de un servidor distribuidor-publicador por lo que una vez acabada la configuración del distribuidor se configura la publicación en el mismo, luego en el servidor de subscripción se configura una subscripción para la publicación anterior.

Una vez creada la instantánea inicial por el agente de la publicación (SnapshotAgent) se trata de establecer una sincronización por el agente de mezcla de la subscripción (MergeAgent). Luego de realizados estos procesos se replica lo filtrado por el publicador hacia el sitio donde defina el suscriptor.

A continuación se explican los procedimientos almacenados fundamentales utilizados:

Distribuidor:

sp_adddistributor (Transact-SQL): configura el servidor como un distribuidor al que se le pasa como parámetro el nombre del servidor de distribución. En este caso se configuró como el nombre del sitio y el nombre de la instancia del servidor (por ejemplo: NombreDePCServidor\NombreDeInstanciaSQLServer). Los otros parámetros que son opcionales se dejaron con el valor por defecto.

`sp_adddistributiondb` (Transact-SQL): crea una nueva BD de distribución e instala el esquema del distribuidor. La BD de distribución almacena procedimientos, esquemas y metadatos usados en la replicación. Este procedimiento almacenado se ejecuta en el distribuidor sobre la BD master en orden para crear la BD de distribución, e instala las tablas necesarias y procedimientos almacenados requeridos para habilitar la distribución de la replicación. Se le pasa como parámetro el nombre de la BD de distribución. En este caso se configuró con el nombre estático “distributor”, entre otros parámetros que son opcionales y se dejaron con el valor por defecto, como por ejemplo el `@security_mode` que tiene valor 1 que significa que el modo de autenticación usado va a ser el de Windows.

`sp_adddistpublisher` (Transact-SQL): configura un publicador para usar una BD de distribución específica. Este procedimiento almacenado se ejecuta en el distribuidor sobre cualquier BD. Es importante señalar que los procedimientos almacenados `sp_adddistributor` (Transact-SQL) y `sp_adddistributiondb` (Transact-SQL) tienen que ser ejecutados antes de éste; se le pasan como parámetros el nombre del servidor de publicación (Publisher) que en este caso coincidirá con el servidor de distribución ya que se está configurando un distribuidor-publicador; además del nombre de la BD de distribución, el `@security_mode` se mantiene en 1 debido al mismo criterio anterior para la autenticación. Es importante señalar que el `@working_directory` con que se va a trabajar es el nombre de la computadora seguido del directorio repldata (por ejemplo: `\\NombreDePCServidor\repldata`).

Publicación:

`sp_replicationdboption` (Transact-SQL): pone las opciones de la BD de replicación para la BD especificada. Este procedimiento almacenado se ejecuta en el publicador o suscriptor en cualquier BD; se le pasan como parámetros: `@dbname` que es el nombre de la BD de la que se va a publicar, `@optname` que se le especifica la opción `mergepublish` lo cual significa que la BD puede ser usada para publicaciones de mezcla. Por último se pone `@value=true` porque si el valor se pone falso y además el `@optname= mergepublish` significa que las subscripciones de las BD de publicación de mezcla son además eliminadas.

`sp_addmergepublication` (Transact-SQL): crea una nueva publicación de mezcla. Este procedimiento almacenado es ejecutado por el publicador en la BD que está siendo publicada; se le pasan como parámetros el nombre de la publicación, el resto se puede dejar con el valor predefinido, aclarando que al menos el `@allow_pull=true` debido a que la subscripción que se va a utilizar es de tipo pull.

`sp_addpublication_snapshot` (Transact-SQL): crea el agente de instantánea para una publicación; se le pasan como parámetros el nombre de la publicación y otro conjunto de entrada que se dejan con los valores predefinidos para cada uno.

`sp_addmergearticle` (Transact-SQL): añade un artículo a una publicación de mezcla existente. Este procedimiento almacenado es ejecutado por el publicador en la BD de publicación; se le pasan como parámetros el nombre de la publicación, nombre del artículo `@article`, objeto de la BD que va a ser publicada `@source_object`, entre otros que se mantuvieron con los valores predefinidos.

Subscripción:

`sp_addmergesubscription` (Transact-SQL): crea una subscripción de mezcla de tipo push o pull. Este procedimiento almacenado se ejecuta por el publicador en la BD de publicación; se le pasan como parámetros el nombre de la publicación, el `@subscriber` que es el nombre del servidor de subscripción, `@subscriber_db` que es el nombre de la BD de subscripción, `@subscription_type` que es el tipo de subscripción a utilizar el cual será pull, entre otros que se mantendrán con los valores predefinidos.

`sp_addmergepullsubscription` (Transact-SQL): agrega una subscripción pull a una publicación de mezcla. Este procedimiento almacenado se ejecuta en el suscriptor sobre la BD de subscripción; se le pasan como parámetros `@publication` que es el nombre de la publicación, `@publisher` que es el nombre del publicador, `@publisher_db` que es el nombre de la BD del publicador y otros con sus valores predefinidos.

`sp_addmergepullsubscription_agent` (Transact-SQL): adiciona un nuevo agente de trabajo usando el esquema de sincronización de subscripciones pull para publicaciones de mezcla. Este procedimiento se ejecuta por el suscriptor sobre la BD de subscripción; se le pasan como parámetros `@publisher` que es el nombre del publicador,

@publisher_db que es el nombre de la BD de publicación, @publication que es el nombre de la publicación, @distributor que es el nombre del distribuidor, entre otros.

De esta manera queda creado el script para cada sitio. Una vez visualizados en la ventana de la aplicación DISTRIBUTOR, se pueden ejecutar en Microsoft SQL Server y efectuar así la replicación.

2.4 Conclusiones parciales

En este capítulo se describieron algunos elementos relacionados con la concepción de la herramienta DISTRIBUTOR; puntualizándose en las partes del diseño de mayor interés como son los diagramas de clases y de secuencia; así como el procedimiento implementado para generación de scripts. En el acápite de Implementación se explicaron todas las clases del proyecto, así como la función de cada uno de sus métodos.

Tomando en consideración los resultados derivados del seguimiento de estos pasos se logró el diseño e implementación de un asistente para la creación de los fragmentos físicos de BDD para el SGBD Microsoft SQL Server.

Capítulo 3. Herramienta para generar esquemas físicos distribuidos en Microsoft SQL Server

En este capítulo se muestra una guía para el uso de la herramienta DISTRIBUTOR. Se enfatiza en el uso de los diferentes elementos presentes en la interfaz gráfica de usuario.

3.1 Presentación de DISTRIBUTOR

Este software cuenta con una ventana principal y otras secundarias. Es completamente operable como un editor de texto que posee varios menús como Archivo, Editar y Ayuda.

El menú Archivo brinda las opciones de:

- Abrir archivo XML (Ctrl+A): Abre un archivo XML.
- Guardar Script (Ctrl+G): Guarda en una carpeta seleccionada varios archivos en texto plano, uno por cada pestaña generada en el software.
- Generar Script (Ctrl+S): Genera el script SQL que posteriormente va a ser utilizado en el Microsoft SQL Server.
- Salir: Cierra o termina la aplicación.

El menú Editar que muestra las opciones:

- Deshacer (Ctrl+Z): Revierte el último cambio realizado en el script.
- Rehacer (Ctrl+Y): Rehace el último cambio realizado.
- Cortar (Ctrl+X): Corta el texto seleccionado en el editor.
- Copiar (Ctrl+C): Copia el texto seleccionado en el editor.
- Pegar (Ctrl+V): Pega el texto anteriormente cortado o copiado en el lugar seleccionado.
- Seleccionar Todo (Ctrl+E): Selecciona todo lo que está en la pestaña.

El menú Ayuda ofrece la opción Acerca de...para mostrar información acerca de DISTRIBUTOR.

En la figura 3.1 se muestra la ventana principal del software DISTRIBUTOR.

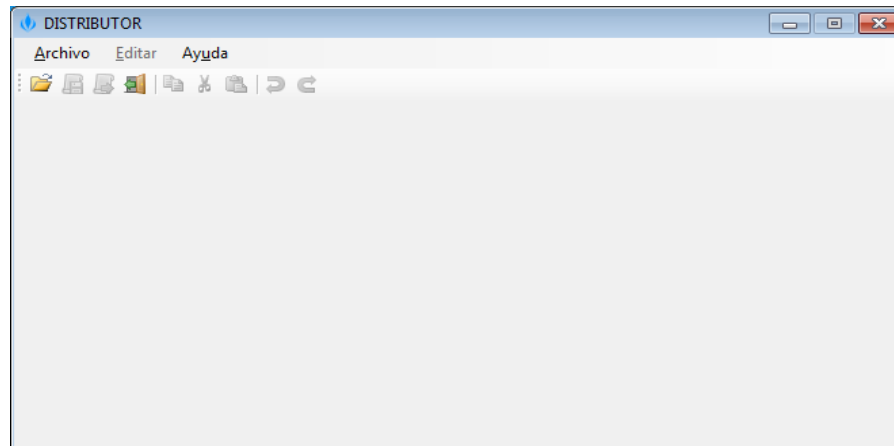


Figura 3.1. Ventana principal de la herramienta DISTRIBUTOR.

Para trabajar con un archivo XML se debe abrir con la opción Abrir archivo XML o mediante las teclas de acceso directo Ctrl+A. Este mostrará un diálogo con el nombre Abrir archivo XML que tiene filtrado para archivos XML (véase la figura 3.2).

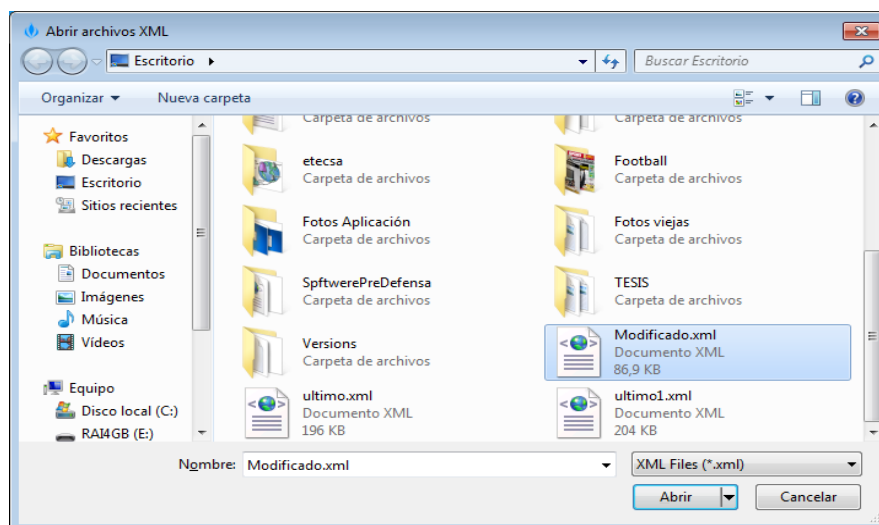


Figura 3.2. Ventana de diálogo para abrir archivo XML.

Una vez cargado el archivo, se generan los nombres de los sitios leídos del XML como pestañas como se muestra en la figura 3.3.

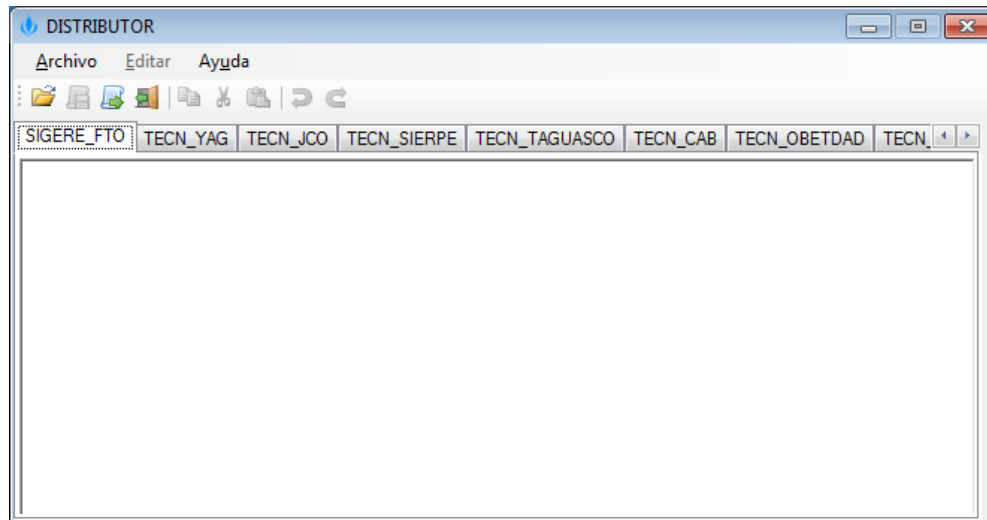


Figura 3.3. Sitios cargados.

En este momento se puede pasar a la opción Generar Script del menú Archivo o utilizar las teclas de acceso directo Ctrl+S.

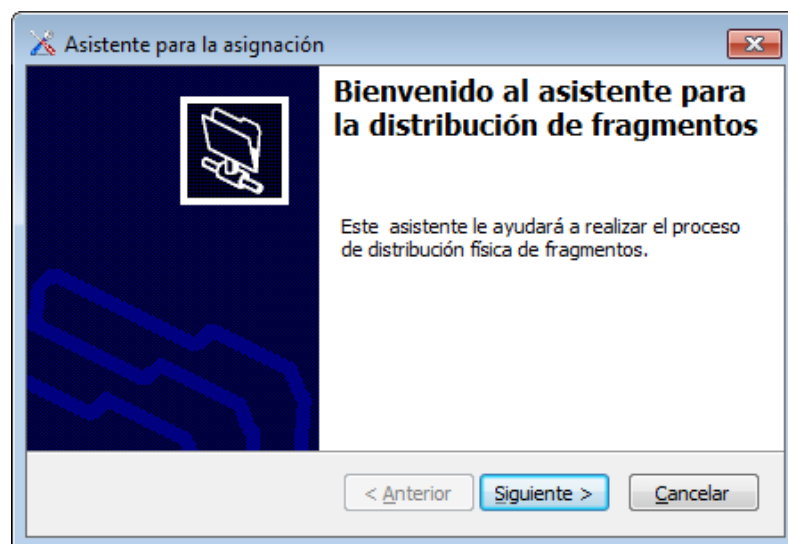


Figura 3.4. Asistente para la asignación.

Al presionar el botón Siguiente o el acceso directo Alt+S, en caso de que haya duplicación de tablas se mostrará un árbol cuyas raíces representan las tablas repetidas y sus hijos los sitios en los que aparecen las mismas; el usuario debe seleccionar en cuál de ellos se encuentra la copia primaria. Obsérvese que en cada ventana se tiene la opción de retroceder un paso atrás, ir hacia adelante o cancelar la operación (véase la figura 3.5).



Figura 3.5. Duplicación.

Una vez que el usuario seleccione dónde se encuentran las copias primarias debe pasar a la siguiente ventana en la cual se obtiene la información para la configuración del servidor de publicación y suscripción, así como la BD de publicación y suscripción. La misma consta de cuatro columnas, la primera se llena automáticamente con los sitios que aparecen en el archivo XML, la segunda corresponde a los servidores de replicación los cuales se buscan de forma automática por la red, la tercera y cuarta columna corresponden a las BD de publicación y suscripción respectivamente (véase la figura 3.6 y 3.7). Es necesario que los servidores de replicación estén disponibles, al menos en este momento, para poder completar el diseño. Esta operación no es restrictiva en modo alguno puesto que cuando la BDD

diseñada se encuentre en funcionamiento, deberá existir conectividad en algún momento.

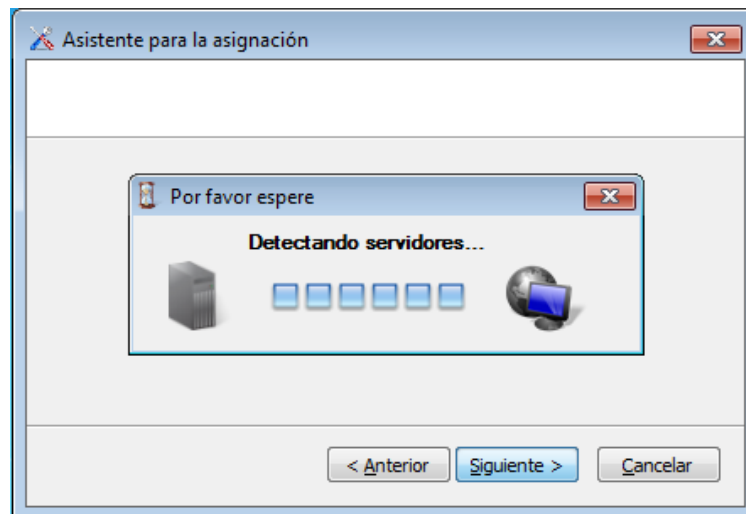


Figura 3.6. Detectando servidores MSQL Server.

Una vez cargados los servidores se muestra la ventana de la figura 3.7.

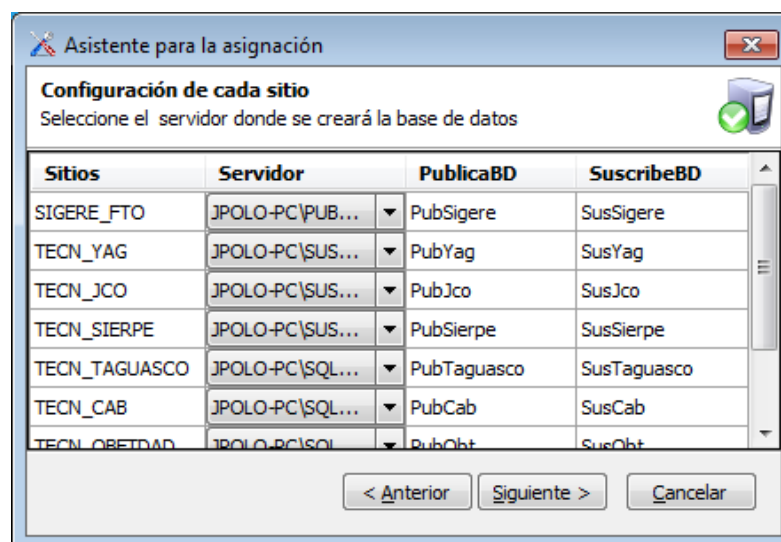


Figura 3.7. Configuración de cada sitio.

En la figura 3.8 se observa la ventana final de ejecución del asistente, culminando así el proceso de distribución física de fragmentos.



Figura 3.8. Fin del asistente.

Al presionar el botón Finalizar, DISTRIBUTOR crea los scripts para cada sitio; dentro de estos se crean las bases de datos, las tablas, se configura la integridad referencial, así como las configuraciones para la replicación (véase la figura 3.9).

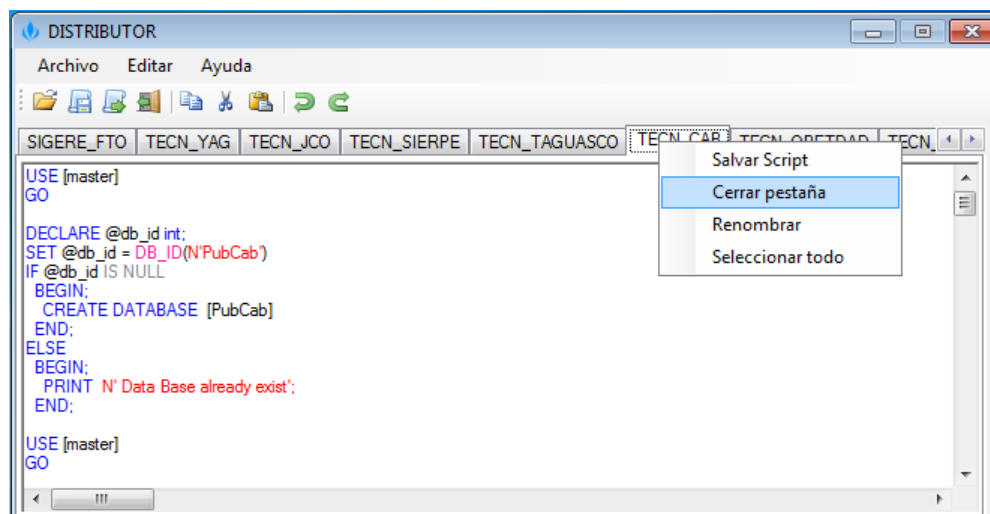


Figura 3.9. Fin de la generación de scripts.

En cada pestaña (TabPage) de la ventana se puede desplegar un menú contextual que brinda las opciones de renombrar Tab, eliminar Tab, seleccionar todo el texto correspondiente al Tab marcado, así como guardar la información del mismo en un archivo de texto txt.

3.2 Descripción de un caso de estudio

La Organización Básica Eléctrica (OBE) de Sancti Spíritus se dio a la tarea de implementar un módulo para el control de los transformadores instalados y las operaciones relacionadas en la explotación de los mismos, reflejando la estructura inherentemente distribuida de las diferentes unidades de la Empresa Eléctrica, teniendo autonomía local y ofreciendo buen rendimiento.

En Cuba, durante los últimos años, la distribución ha sido el área de trabajo menos atendida de la Unión Eléctrica, presentándose altos índices de pérdidas e interrupciones, un creciente número de transformadores dañados, y un escaso nivel de informatización y automatización.

Como respuesta a los problemas informáticos detectados, se identifica como necesidad el desarrollo de un sistema que permita controlar los transformadores instalados con una reducción de los gastos totales de explotación de la distribución, influyendo en un mejor servicio a los clientes. El transformador es el equipo más cercano al cliente que más abunda en las redes eléctricas cubanas, con una distribución espacial muy variada y el mayor índice de fallas; de ahí la importancia de minimizar sus averías. Esta es otra razón por la que se necesita desarrollar un sistema capaz de automatizar toda la información relacionada con ellos, poder seguir sus ciclos de mantenimiento, su estado de carga; así como ser capaz de prevenir las fallas.

Para diseñar el sistema hay que tener en cuenta la estructura geográficamente distribuida de las Empresas Eléctricas y sus propias características de comunicaciones; reconociendo los cuatro niveles bien definidos: nacional, provincial, territorial y de sucursal, además de un quinto nivel que es el del taller, al mismo nivel de la provincia.

La comunicación entre la OBE Provincial y las OBEs territoriales se hace a través de líneas telefónicas de baja velocidad; aunque en los últimos tiempos se han introducido gradualmente las redes inalámbricas, aún no se cuenta con este sistema en todas las provincias, ni es política de la Unión Eléctrica hacerlo extensivo en poco tiempo. Como la topología de las redes informáticas no es la más adecuada para contar con una BDC provincial donde accedan todos los usuarios, ya sean de la provincia, taller o territorio, es necesario diseñar una BDD, para poder seguir el ciclo de vida completo de cada transformador manteniendo datos históricos de todas las operaciones que se realizan.

Esto le da la facilidad de hacer numerosos estudios comparativos del comportamiento de las diferentes instalaciones, equipos, así como prever posibles fallas en ellos. El transformador tiene un ciclo de vida muy difícil de seguir ya que pasa por diferentes estados en diferentes ubicaciones.

El control de los transformadores se hace en el ámbito provincial, por lo que sólo intervienen la provincia, el territorio y el taller como sitios del sistema. El ciclo de vida del transformador comienza en el taller mediante la verificación de un transformador nuevo.

Después pasa a formar parte de los transformadores disponibles en la OBE Provincial donde aparecerán todos los datos de su ficha. Allí, dependiendo de sus características y las necesidades de las diferentes OBEs territoriales, es asignado a una de estas, y de ahí hacia el banco que lo necesita. En el municipio sólo se encuentran los datos pertenecientes a las instalaciones y equipos que atiende este territorio, mientras que en la provincia está centralizada toda la información de las OBEs territoriales.

En la OBE territorial se realizan todas las operaciones sobre el transformador que van conformando una historia. Existen dos formas para que un transformador retorne hacia la provincia: moviéndolo directo al taller o mediante un evento de Reporte de inspección al transformador dañado que lo retira automáticamente del Banco y lo envía al taller si se comprueba que está dañado, a la vez que se crea un evento de Necesidad de transformador para sustituir este averiado. Luego de llegar el transformador y sus datos

al taller, se realiza la defectación o diagnóstico en el Taller. Esta operación es la que dice si el transformador se retira definitivamente por su daño o si su problema es solucionable. En este caso, el transformador no pierde su historia, se da como disponible, pasa a la provincia y vuelve a fluir, dependiendo de sus características técnicas y las necesidades de las OBEs territoriales. Nótese que en todo este ir y venir no pierde los datos de todas las pruebas que se le hayan realizado, aunque sea situado en una OBE territorial diferente.

Para este caso de estudio se tienen los siguientes sitios, donde para cada uno de ellos se va a generar un script mediante DISTRIBUTOR. Nótese que todo el diseño de distribución precedente a la generación física de esquemas es realizado por las herramientas integradas en SIADBDD cuyos resultados quedan depositados en un archivo XML que sirve de entrada a DISTRIBUTOR (véase Anexo 1).

- UEB Subcentro Fomento (Sigere_Fto)
- UEB OBE Sancti Spíritus (Tecn_ssp)
- UEB Subcentro La Sierpe (Tecn_Sierpe)
- UEB Subcentro Taguasco (Tecn_Taguasco)
- UEB OBE Trinidad (Tecn_Obetdad)
- UEB OBE Yaguajay (Tecn_Yag)
- UEB OBE Cabaiguán (Tecn_Cab)
- UEB OBE Jatibonico (Tecn_Jco)
- Taller (Transformadores)

3.3 Solución

Para utilizar en Microsoft SQL Server los resultados obtenidos por DISTRIBUTOR se debe abrir el administrador Microsoft SQL Server Management Studio. Posteriormente se copia el script hacia una nueva consulta y se ejecuta (véase Anexo 2). Una vez realizada esta acción se crean las BD así como las tablas referentes a ellas; se configura los servidores de distribución, publicación y suscripción, además de las publicaciones y suscripciones que contiene el script para ese sitio, quedando como se muestra en la figura 3.10.

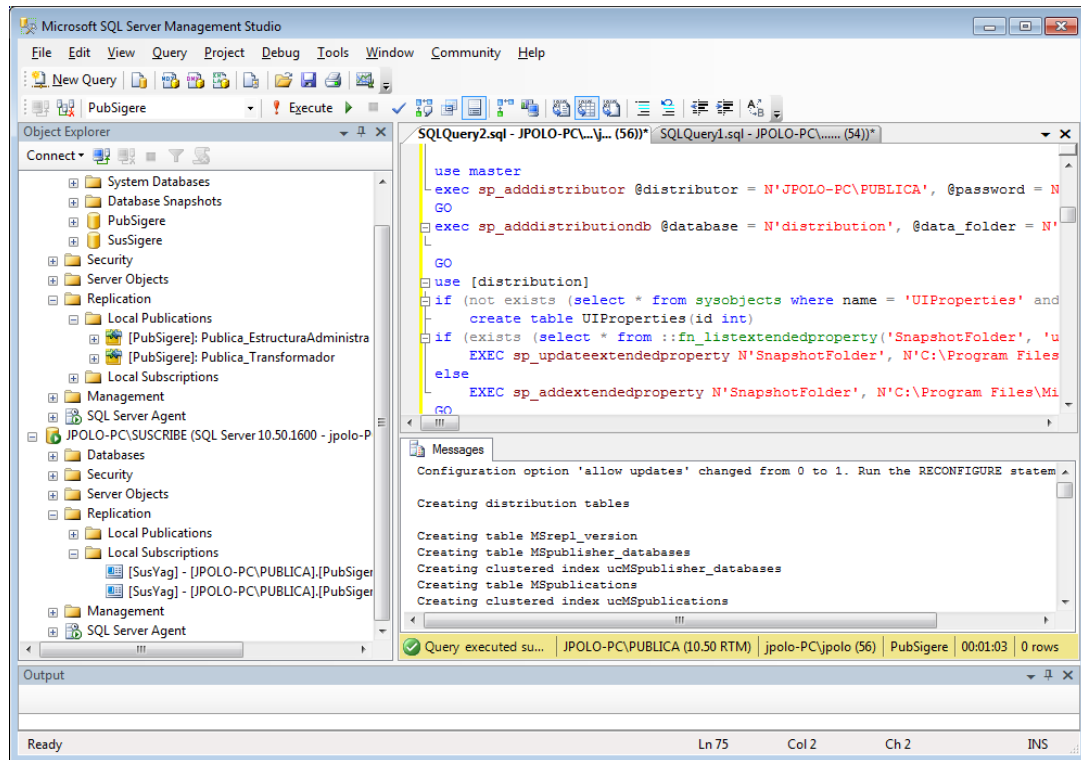


Figura 3.10. Generación de scripts.

Luego de ser ejecutado el script, se crea la instantánea inicial por el agente de instantáneas (SnapshotAgent). Para ello se debe hacer clic derecho sobre el nombre de la publicación ([PubSigere]: Publica_EstructuraAdministra) y elegir la opción View SnapshotAgent Status para ver el estado.

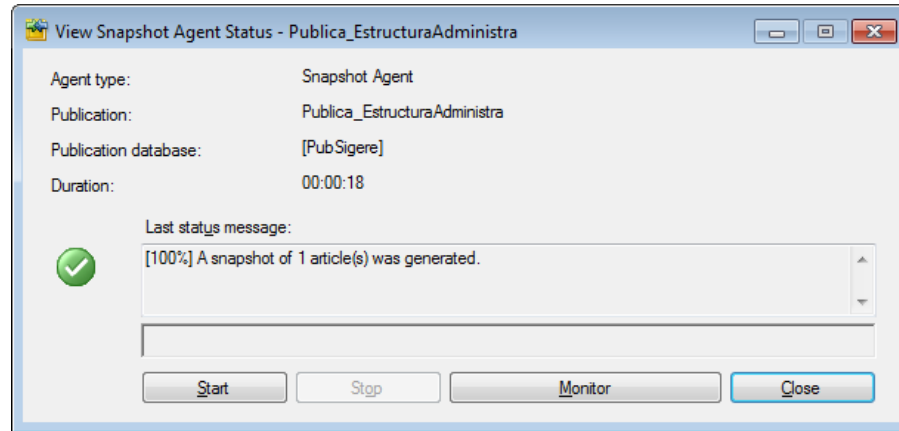


Figura 3.11. Agente de instantáneas.

Seguidamente se puede comprobar que fue realizada la sincronización haciendo clic derecho sobre la subscripción ([SubYag]-[JPOLO-PC\PUBLICA].[PubSigere]: Publica_EstructuraAdministra) y elegir la opción View Synchronization Status, el cual brinda la información que se observa en la figura 3.12.

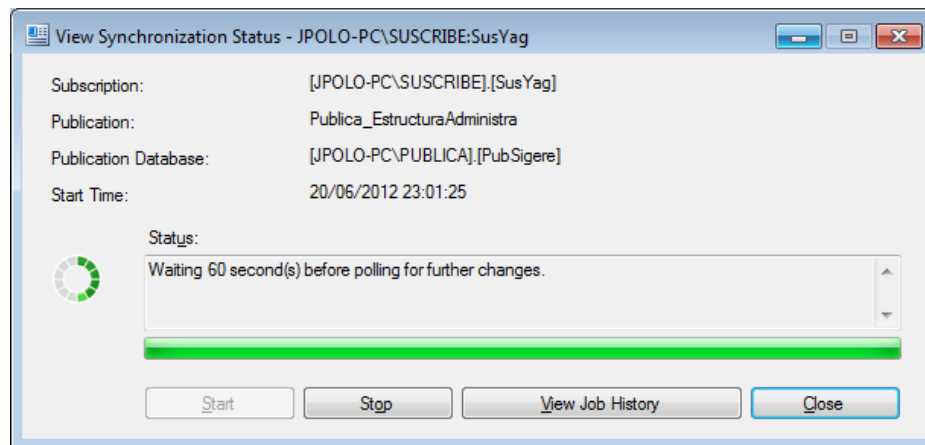


Figura 3.12. Estado de sincronización.

Una vez comprobadas las acciones anteriores se puede probar que se haya efectuado la replicación, para lo que se debe expandir la base de datos de subscripción en el panel izquierdo, con la finalidad de chequear si se escribieron los cambios en la misma; en este caso se ejecutó perfectamente el modelo de réplica, las 16 columnas que posee

la tabla EstructuraAdministra se escribieron en la base de datos de subscripción SubYag (véase la figura 3.13).

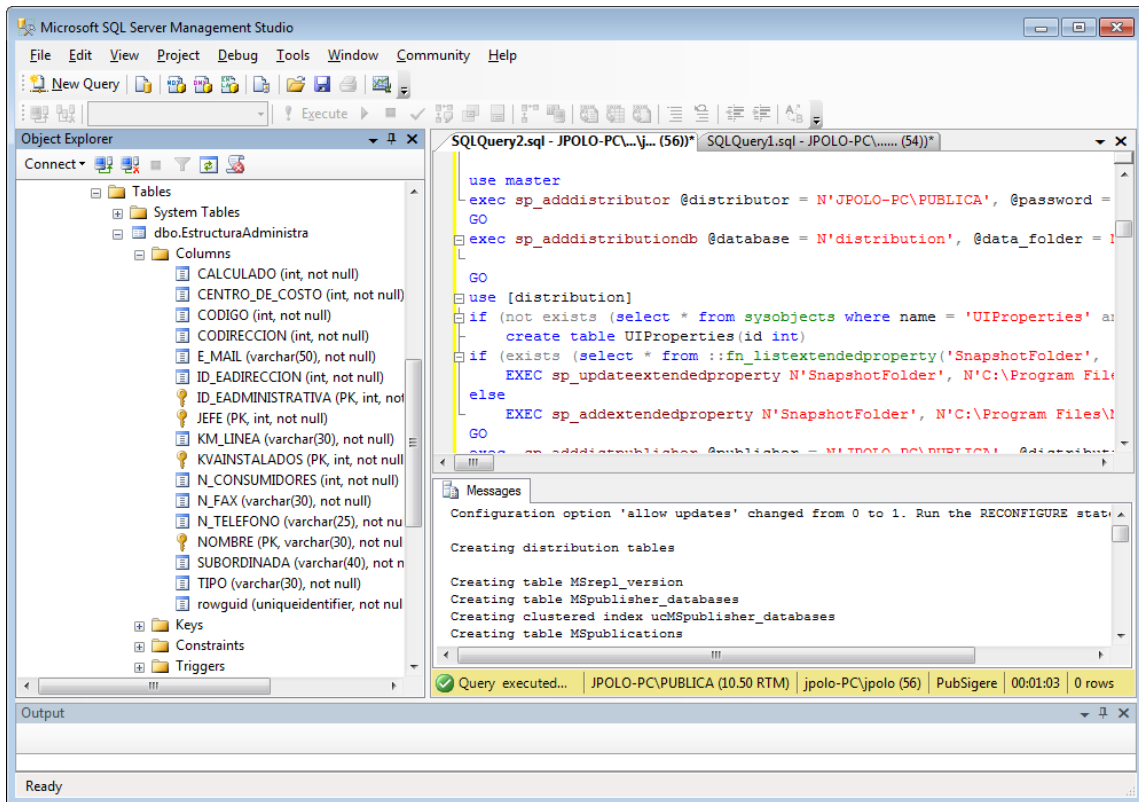


Figura 3.13. Estado de replicación.

De esta manera se ven los resultados del filtrado aplicado a la tabla de publicación.

3.4 Consideraciones parciales

En este capítulo se muestra cómo trabajar con la herramienta DISTRIBUTOR y se detallan aspectos de la prueba realizada con un caso de estudio real, tomando como punto de partida el archivo XML generado para este caso sobre transformadores de distribución de energía eléctrica. A partir de este archivo se obtuvo un conjunto de scripts para la generación de esquemas físicos en cada sitio. Estos fueron ejecutados en Microsoft SQL Server lográndose la replicación de datos.

Conclusiones

1. Se determinaron elementos de configuración de los servidores, tipo de replicación, publicación y suscripción apropiados para el diseño de BDD en Microsoft SQL Server.
2. Se identificaron los aspectos de programación en Transact-SQL mediante llamada a procedimientos almacenados para configurar un entorno de replicación completo que permitió materializar diseños de BDD a nivel físico en Microsoft SQL Server.
3. Se implementó una nueva versión de la herramienta DISTRIBUTOR que sigue el diseño realizado, es fácil de usar y representa una solución al problema de ubicación que es independiente de las estructuras internas manejadas por SIADBDD y permite la configuración avanzada de parámetros relativos al entorno de replicación en Microsoft SQL Server.
4. Se dio mantenimiento correctivo y perfectivo a la herramienta DISTRIBUTOR que permite generar los esquemas físicos de ubicación usando los elementos programáticos de replicación identificados.
5. Se probó la herramienta con un caso de estudio real.

Recomendaciones

1. La implementación de salidas para otros gestores de bases de datos como Postgresql y Oracle.
2. Incluir en el Asistente para la asignación la opción de elegir de forma dinámica cuál de las 3 tipologías de replicación se va a tener en cuenta a la hora de generar los scripts.

Referencias Bibliográficas

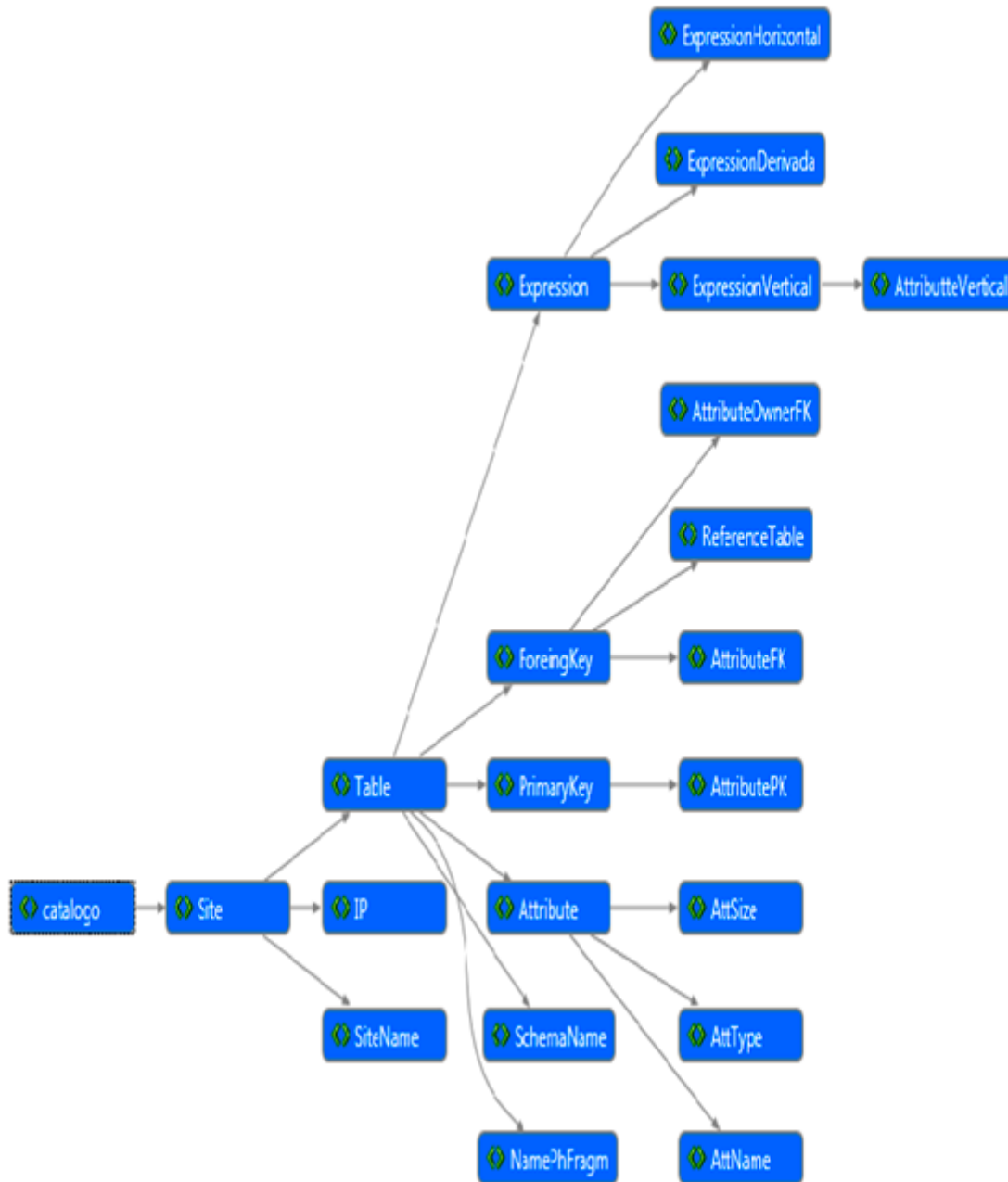
- BAIÃO, F. A., MATTOSO, M., SHAVLIK, J. W. & ZAVERUCHA, G. Applying Theory Revision to the Design of Distributed Databases. *In: HORVÁTH, T. & YAMAMOTO, A., eds. Proceedings of the 13th International Conference on Inductive Logic Programming ILP 2003, LNAI 2835, September 29-October 1 2003 Szeged, Hungary. Springer, 57-74.*
- BAIÃO, F. A., MATTOSO, M. & ZAVERUCHA, G. A Framework for the Design of Distributed Databases. *In: LITWIN, W. & LÉVY, G., eds. Records of the 4th International Meeting on Distributed Data & Structures 4 (WDAS 2002), 2002 Paris, France. Carleton Scientific, 29-36.*
- BAIÃO, F. A., MATTOSO, M. & ZAVERUCHA, G. 2004. A Distribution Design Methodology for Object DBMS. *Distributed and Parallel Databases*, 16, 45-90.
- BELLATRECHE, L., KARLAPEM, K. & SIMONET, A. 2000. Algorithms and support for horizontal class partitioning in object-oriented databases. *Distributed and Parallel Databases*, 8, 155-179.
- BERTONE, R. 2004. Métricas de Performance en Administración de BD Distribuidas en redes LAN y WAN.
- BHALLA, S. & HASEGAWA, M. Parallelizing serializable transactions within distributed real-time database systems Proceedings of the International Conference on Embedded and Ubiquitous Computing (EUC 2005). Lecture Notes in Computer Science, 2005. Springer, 203-213
- CERI, S., MARTELLA, G. & PELAGATTI, G. 1982. Optimal File Allocation in a Computer Network: a Solution Method Based on the Knapsack Problem. *Computer Networks*, 6, 345-357.
- CERI, S., NAVATHE, S. B. & WIEDERHOLD, G. 1983. Distribution Design of Logical Database Schemas. *IEEE Trans. Software Eng.*, 9, 487-504.
- CERI, S. & PERNICI, B. 1985. *DATAID-D: Methodology for Distributed Database Design*, North-Holland.
- CERI, S., PERNICI, B. & WIEDERHOLD, G. 1987. Distributed Database Design Methodologies. *IEEE Database Eng. Bull.*, 75, 533-546.
- COULON, C., PACITTI, E. & VALDURIEZ, P. Consistency Management for Partial Replication in a High Performance Database Cluster. IEEE International Conference on Parallel and Distributed Systems (ICPADS2005), 2005 Fukuoka, Japan.
- DATE, C. J. 2000. *Introducción a los Sistemas de Bases de Datos, Séptima edición*, México, Addison-Wesley.
- GANÇARSKI, S., NAACKE, H., PACITTI, E. & VALDURIEZ, P. Parallel Processing with Autonomous Databases in a Cluster System. *In: MEERSMAN, R. & TARI, Z., eds. Proceedings of the Confederated International Conferences On the Move to Meaningful Internet Systems, 2002 - DOA/CoopIS/ODBASE, October 30-November 1, 2002 Irvine, California, USA. Springer, 410-428.*
-

- HABABEH, I. O., BOWRING, N. & RAMACHANDRAN, M. A Method for Fragment Allocation in Distributed Object Oriented Database Systems. Proceedings of the 5th Annual PostGraduate Symposium on The Convergence of Telecommunications, Networking & Broadcasting (PGNet), June, 2004 2004 Liverpool, UK. 54 - 59.
- JOHANSSON, J. M., MARCH, S. T. & NAUMANN, J. D. The effects of parallel processing on update response time in distributed database design. *In*: ANG, S., KRCCMAR, H., ORLIKOWSKI, W. J., WEILL, P. & DEGROSS, J. I., eds. Proceedings of the Twenty-First International Conference on Information Systems ICIS, December 10-13, 2000 2000 Brisbane, Australia. ACM, 187-196.
- LEE, J. W. & BAIK, D. K. 2004. Database allocation modeling for optimal design of distributed systems. *IEICE Transactions on Information and Systems*, E87D, 1795-1804.
- LIN, Y., KEMME, B., PATIÑO-MARTÍNEZ, M. & JIMÉNEZ-PERIS, R. Middleware based Data Replication providing Snapshot Isolation. Proceedings of the ACM SIGMOD International Conference on Management of Data., June 14-16, 2005 2005 Baltimore, Maryland, USA. ACM Press, 419-430.
- MA, H., SCHEWE, K.-D. & WANG, Q. 2005. Distribution design for higher-order data models. Palmerston North, New Zealand: Massey University, Department of Information Systems.
- MEI, H. & SHENG, O. A Semantic Based Methodology for Integrated Computer-Aided Distributed Database Design. The 25th Hawaii International Conference on System Sciences, 1992. 288- 299.
- MORELL, D. E. C. 2004. Replicación de Datos en SQL Server.
- NAVATHE, S. B., KARLPALEM, K. & RA, M. 1995. A mixed fragmentation methodology for initial distributed database design. Atlanta, GA: College of Computing, Georgia Institute of Technology.
- ÖZSU, M. T. & VALDURIEZ, P. 1991. *Principles of Distributed Database Systems*, Prentice-Hall.
- ÖZSU, M. T. & VALDURIEZ, P. 1999. *Principles of Distributed Database Systems, Second Edition*, Upper Saddle River, New Jersey., Prentice-Hall.
- PÉREZ, J., PAZOS, R. A., FRAUSTO-SOLIS, J., REYES, G., SANTAOLAYA, R., FRAIRE, H. J. & CRUZ, L. An approach for solving very large scale instances of the design distribution problem for distributed database systems. Proceedings of the 4th International School and Symposium on Advanced Distributed Systems (ISSADS2005). Lecture Notes in Computer Science, 2005. Springer, 33-42.
- SAVONNET, M., TERRASSE, M.-N. & YÉTONGNON, K. FRAGTIQUE: An OO Distribution Design Methodology. *In*: CHEN, A. L. P. & LOCHOVSKY, F. H., eds. Proceedings of the Sixth International Conference on Database Systems for Advanced Applications (DASFAA), April 19-21 1999 Hsinchu, Taiwan. IEEE Computer Society, 283-290.
- SCHEWE, K.-D. Fragmentation of object oriented and semi-structured data. *In*: HAAV, H.-M. & KALJA, A., eds. Databases and Information Systems II, 2002. Kluwer Academic Publishers, 1-14.

TAMHANKAR, A. M. & RAM, S. 1998. Database fragmentation and allocation: an integrated methodology and case study. *IEEE Transactions on Systems, Man, and Cybernetic Part A*, 28, 288-305.

Anexos

Anexo 1. Esquema XML



Anexo 2. Script generado por DISTRIBUTOR

```
USE [master]
GO
```

```
DECLARE @db_id int;
SET @db_id = DB_ID(N'PubSigere')
IF @db_id IS NULL
    BEGIN;
        CREATE DATABASE [PubSigere]
    END;
ELSE
    BEGIN;
        PRINT N' Data Base already exists';
    END;
```

```
USE [master]
GO
```

```
DECLARE @db_id int;
SET @db_id = DB_ID(N'SusSigere')
IF @db_id IS NULL
    BEGIN;
        CREATE DATABASE [SusSigere]
    END;
ELSE
    BEGIN;
        PRINT N' Data Base already exists';
    END;
```

```
use master
exec sp_adddistributor @distributor = N'JPOLO-PC\PUBLICA', @password = N''
GO
exec sp_adddistributiondb @database = N'distribution', @data_folder = N'C:\Program
Files\Microsoft SQL Server\MSSQL10_50.PUBLICA\MSSQL\Data', @log_folder =
N'C:\Program Files\Microsoft SQL Server\MSSQL10_50.PUBLICA\MSSQL\Data',
@log_file_size = 2, @min_distretention = 0, @max_distretention = 72,
@history_retention = 48, @security_mode = 1
```

```
GO
```

```

use [distribution]
if (not exists (select * from sysobjects where name = 'UIProperties' and type = 'U '))
    create table UIProperties(id int)
if (exists (select * from ::fn_listextendedproperty('SnapshotFolder', 'user', 'dbo', 'table',
'UIProperties', null, null)))
    EXEC sp_updateextendedproperty N'SnapshotFolder', N'C:\Program
Files\Microsoft SQL Server\MSSQL10_50.PUBLICA\MSSQL\RepIData', 'user', dbo,
'table', 'UIProperties'
else
    EXEC sp_addextendedproperty N'SnapshotFolder', N'C:\Program Files\Microsoft
SQL Server\MSSQL10_50.PUBLICA\MSSQL\RepIData', 'user', dbo, 'table',
'UIProperties'
GO
exec sp_adddistpublisher @publisher = N'JPOLO-PC\PUBLICA', @distribution_db =
N'distribution', @security_mode = 1, @working_directory = N'C:\Program Files\Microsoft
SQL Server\MSSQL10_50.PUBLICA\MSSQL\RepIData', @trusted = N'false',
@thirdparty_flag = 0, @publisher_type = N'MSSQLSERVER'
GO

use [PubSigere]
exec sp_replicationdboption @dbname = N'PubSigere', @optname = N'merge publish',
@value = N'true'
GO
use[PubSigere]
CREATE TABLE      EstructuraAdministra
(
    CALCULADO          int NOT NULL,
    CENTRO_DE_COSTO    int NOT NULL,
    CODIGO              int NOT NULL,
    CODIRECCION         int NOT NULL,
    E_MAIL              varchar(50) NOT NULL,
    ID_EADIRECCION      int NOT NULL,
    ID_EADMINISTRATIVA  int NOT NULL,
    JEFE                int NOT NULL,
    KM_LINEA            varchar(30) NOT NULL,
    KVAINSTALADOS       int NOT NULL,
    N_CONSUMIDORES      int NOT NULL,
    N_FAX               varchar(30) NOT NULL,
    N_TELEFONO          varchar(25) NOT NULL,
    NOMBRE              varchar(30) NOT NULL,
    SUBORDINADA         varchar(40) NOT NULL,
    TIPO                varchar(30) NOT NULL,

```

```

CONSTRAINT    PK_EstructuraAdministra PRIMARY KEY CLUSTERED
(
    ID_EADMINISTRATIVA ASC,
    JEFE ASC,
    KVAINSTALADOS ASC,
    NOMBRE ASC
)
)
GO
-- Enabling the replication database
use [PubSigere]
exec sp_replicationdboption @dbname = N'PubSigere', @optname = N'merge publish',
@value = N'true'
GO
-- Adding the merge publication
use [PubSigere]
exec sp_addmergepublication @publication = N'Publica_EstructuraAdministra',
@description = N'Merge publication of database "PubSigere" from Publisher "JPOLO-
PC\PUBLICA".', @sync_mode = N'native', @retention = 14, @allow_push = N'true',
@allow_pull = N'true', @allow_anonymous = N'true', @enabled_for_internet = N'false',
@snapshot_in_defaultfolder = N'true', @compress_snapshot = N'false', @ftp_port = 21,
@ftp_subdirectory = N'ftp', @ftp_login = N'anonymous', @allow_subscription_copy =
N'false', @add_to_active_directory = N'false', @dynamic_filters = N'false',
@conflict_retention = 14, @keep_partition_changes = N'false', @allow_synctoalternate
= N'false', @max_concurrent_merge = 0, @max_concurrent_dynamic_snapshots = 0,
@use_partition_groups = null, @publication_compatibility_level = N'100RTM',
@replicate_ddl = 1, @allow_subscriber_initiated_snapshot = N'false',
@allow_web_synchronization = N'false', @allow_partition_realignment = N'true',
@retention_period_unit = N'days', @conflict_logging = N'both',
@automatic_reinitialization_policy = 0
GO
exec sp_addpublication_snapshot @publication = N'Publica_EstructuraAdministra',
@frequency_type = 4, @frequency_interval = 14, @frequency_relative_interval = 1,
@frequency_recurrence_factor = 0, @frequency_subday = 1,
@frequency_subday_interval = 5, @active_start_time_of_day = 500,
@active_end_time_of_day = 235959, @active_start_date = 0, @active_end_date = 0,
@job_login = null, @job_password = null, @publisher_security_mode = 1
GO
use [PubSigere]
exec sp_addmergearticle @publication = N'Publica_EstructuraAdministra', @article =
N'EstructuraAdministra', @source_owner = N'dbo', @source_object =
N'EstructuraAdministra', @type = N'table', @description = null, @creation_script = null,
@pre_creation_cmd = N'drop', @schema_option = 0x0000000010C034FD1,
@identityrangemanagementoption = N'manual', @destination_owner = N'dbo',
@force_reinit_subscription = 1, @column_tracking = N'false', @subset_filterclause =

```

```

null, @vertical_partition = N'false', @verify_resolver_signature = 1,
@allow_interactive_resolver = N'false', @fast_multicol_updateproc = N'true',
@check_permissions = 0, @subscriber_upload_options = 0, @delete_tracking = N'true',
@compensate_for_errors = N'false', @stream_blob_columns = N'false',
@partition_options = 0
GO
use [PubSigere]
exec sp_addmergesubscription @publication = N'Publica_EstructuraAdministra',
@subscriber = N'JPOLO-PC\SUSCRIBE', @subscriber_db = N'SusYag',
@subscription_type = N'pull', @subscriber_type = N'local', @subscription_priority = 0,
@sync_type = N'Automatic'

exec sp_addmergesubscription @publication = N'Publica_EstructuraAdministra',
@subscriber = N'JPOLO-PC\SUSCRIBEOTRA', @subscriber_db = N'SusJco',
@subscription_type = N'pull', @subscriber_type = N'local', @subscription_priority = 0,
@sync_type = N'Automatic'
GO
exec sp_startpublication_snapshot @publication = N'Publica_EstructuraAdministra'
GO
-----BEGIN: Script to be run at Subscriber "-----

use [SusSigere]
exec sp_addmergepullsubscription @publisher = N'JPOLO-PC\SUSCRIBE',
@publication = N'Publica_Transformador', @publisher_db = N'PubYag',
@subscriber_type = N'Local', @subscription_priority = 0, @description = N'',
@sync_type = N'Automatic'

exec sp_addmergepullsubscription_agent @publisher = N'JPOLO-PC\SUSCRIBE',
@publisher_db = N'PubYag', @publication = N'Publica_Transformador', @distributor =
N'JPOLO-PC\SUSCRIBE', @distributor_security_mode = 1, @distributor_login = N'',
@distributor_password = null, @enabled_for_syncmgr = N'False', @frequency_type =
64, @frequency_interval = 0, @frequency_relative_interval = 0,
@frequency_recurrence_factor = 0, @frequency_subday = 0,
@frequency_subday_interval = 0, @active_start_time_of_day = 0,
@active_end_time_of_day = 235959, @active_start_date = 20120529,
@active_end_date = 99991231, @alt_snapshot_folder = N'', @working_directory = N'',
@use_ftp = N'False', @job_login = null, @job_password = null,
@publisher_security_mode = 1, @publisher_login = null, @publisher_password = null,
@use_interactive_resolver = N'False', @dynamic_snapshot_location = null,
@use_web_sync = 0
GO
-----END: Script to be run at Subscriber "-----

use[PubSigere]

```

```

CREATE TABLE    OTRATABLA
(
    CIRCUITO    varchar(20) NOT NULL,
    CODIGO      int NOT NULL,
    CODIGOANTIGUO int NOT NULL,
    CODIRECCION int NOT NULL,
    CONEXION    varchar(50) NOT NULL,
    ESTADOOOPERATIVO    varchar(30) NOT NULL,
    ID_EADIRECCION int NOT NULL,
    ID_EADMINISTRATIVA int NOT NULL,
    ID_VOLTAJEPRIMARIO int NOT NULL,
    ID_VOLTAJESALIDA    int NOT NULL,
    NSECCION    varchar(50) NOT NULL,
    NUMCLIENTES int NOT NULL,
    SECCIONALIZADOR    varchar(50) NOT NULL,
    TIPOALIMENTACION    varchar(30) NOT NULL,
    TIPOSALIDA varchar(50) NOT NULL,
    TIPOTERRENO    varchar(30) NOT NULL,
    CONSTRAINT    PK_OTRATABLA    PRIMARY KEY CLUSTERED
        (
            CODIGO    ASC
        )
)
GO
use [SusSigere]
exec sp_resyncmergesubscription @publisher =N'JPOLO-PC\SUSCRIBE',
    @publisher_db = N'PubYag', @publication = N'Publica_Transformador', @resync_type
= 0
GO

```