

Notificaciones Push

Dispositivos móviles y web

Guillermo López Leal

Notificaciones Push: Dispositivos móviles y web

Guillermo López Leal

Copyright © 2012 Guillermo López Leal, Telefónica Digital (PDI), All rights reserved.

Agradecimientos

- A Bernardo López y Carmen Leal, por todo lo que han hecho siempre por mí. Gracias.
- A Leticia Núñez, por aguantar mis peculiaridades.
- A Francisco Marzal, por su inestimable ayuda durante la carrera y el proyecto.
- A Fernando Rodríguez Sela, Fernando Jiménez Moreno, José Antonio Olivera Ortega, Ignacio Eliseo Barandalla Torregrosa, Daniel Coloma, Miguel Schneider y al resto de gente de Telefónica I+D, ya sea de la iniciativa Open Web Device u otras, por todo lo que me han enseñado durante todo este tiempo. Gracias a todos.
- A Manuel Touriño y Francisco Sánchez, por la revisión a fondo de este proyecto, por estar cuando se les necesita y por ser grandes compañeros de carrera.
- A todos mis compañeros de Mozilla Hispano, por todo lo que me han enseñado estos últimos años, y por todo lo que hemos logrado juntos.
- A mis amigos de toda la vida, y a los más recientes, por ser como son y pasar enormes tiempos de risa juntos.

Tabla de contenidos

1. Introducción	1
Resumen	2
Objetivos	2
Descripción del servicio	3
2. Estado del arte	4
Redes móviles: funcionamiento y problemas	4
Redes LAN públicas o privadas	4
Estados del dominio de circuitos	6
Operadores	11
WAP Push	12
Internet	14
GCM: Google Cloud Messaging	15
APNS: Apple Push Notification Service	20
Thiafi: el futuro de Google	22
3. API	27
Externas	28
Entre la WA y UA	28
Entre el WA y AS	33
Entre el AS y NS	33
Interna: Entre el UA y el NS	38
Interfaz HTTP	39
Interfaz WebSocket	39
Interfaz UDP	49
Flujo completo	50
4. Arquitectura interna	51
Tipos de servidor	51
NS-UA-WS	52
NS-UA-UDP	52
NS_Proxy_UDP	53
NS-AS	53
NS-MSG-Monitor	53
Colas de mensajes, o Message Queue	53
Base de datos, NO-SQL (MongoDB)	54
Flujos de trabajo	54
Registros de dispositivos/aplicaciones	54
Aceptación de notificación	55
Envío de notificación al dispositivo: WebSocket	56
Envío de notificación al dispositivo: UDP	56
Confirmación de recepción	57
5. Tecnologías	58
Node.js	58
Para qué se usa	60
Por qué se ha elegido	60
RabbitMQ	61

Para qué se usa	62
Por qué se ha elegido	62
MongoDB	63
Para qué se usa	65
Por qué se ha elegido	65
6. Seguridad	66
Tokens de dispositivo	66
Registro con clave pública-privada	68
Notificaciones firmadas	69
Verificación de notificaciones	70
Ataque DDoS: flooding y replay	71
Comunicación cifrada vía SSL	72
UDP para notificaciones	73
Entrega de notificaciones a las aplicaciones	73
7. Lecciones aprendidas	75
SQL contra NoSQL	75
Investigar nuevas tecnologías	76
Tener tests unitarios y funcionales	77
Priorización y enfoque en tareas	77
Trabajar en abierto	78
La importancia del hardware	78
8. Futuro: v2 y siguientes	79
Sobreescribir notificaciones	79
Recepción segura de token de usuario	80
Cambio en el tamaño máximo	80
Poder enviar números de versión	80
Implementar soporte para IPv6	81
Implementar prioridades	83
PINGs de backup	83
Soporte WAP push	84
Control de abuso	84
Soporte a diferentes servidores de notificaciones	84
Control de presencia	85
Control de entrega de notificaciones	86
9. Conclusiones	87
Glosario de términos	89
Bibliografía	90

Lista de figuras

2.1. Esquema de estados radio	8
2.2. Aplicación con demasiado tráfico de control, 1	9
2.3. Aplicación con demasiado tráfico de control, 2	10
2.4. Aplicación con poco tráfico de control	11
2.5. Funcionamiento WAP Push	13
2.6. GCM logo	16
2.7. Funcionamiento de GCM	19
2.8. Líneas de código para implementar Thialfi	24
3.1. Instancias y comunicaciones	28
3.2. Interfaz <code>PushManager</code>	29
3.3. Función <code>requestRemotePermission</code>	29
3.4. Función <code>revokeRemotePermission</code>	32
3.5. Notificación enviada	34
4.1. Esquema general de la arquitectura	51
4.2. Flujo de registro	55
4.3. Flujo de notificación	55
6.1. Generación de token	67
6.2. Verificación de tokens	68
6.3. Registro de WA en el UA	69
6.4. Verificación de firma (en <code>src/common/crypto.js</code>)	69
6.5. Comprobación de firma (en <code>src/ns_as/ns_as_server.js</code>)	70
6.6. Notificación normalizada	71

Lista de ejemplos

3.1. Pedir permiso remoto y recibir URL	31
3.2. Pedir revocación de permiso (revocar URL)	33
8.1. Notificación sobre el desarrollo de un partido de fútbol	79
8.2. Uso de números de versión en aplicación del banco	81
8.3. Información sobre lugares cercanos y ofertas	83
8.4. Enviar mensajes sólo si los anteriores han llegado	86

Capítulo 1. Introducción

Una de las grandes ventajas que ha aportado Internet al desarrollo de aplicaciones es, aparte de una distribución mucho más sencilla, la posibilidad de poder obtener datos de forma dinámica. La mayoría de las veces se permite cambiar el comportamiento del propio programa gracias a los datos descargados desde la red o incluso actualizar la información que se va a mostrar a los usuarios dependiendo del momento del día u otras variables: los resultados de nuestro equipo favorito de fútbol, el tiempo que va a hacer en nuestra ciudad, o incluso mensajería instantánea, son algunos de los usos que se hacen de las tecnologías de Internet para aplicaciones locales.

A lo largo de los años, las tecnologías para descargar estos datos han ido variando, acomodándose a los cambios que han sucedido en Internet, desde las redes de módem lentas que se usaban hace un tiempo atrás, hasta redes de fibra óptica capaces de descargar un gran flujo de datos pasando por las olvidadas redes móviles, que son aquellas que más han sufrido esta descarga de datos de una manera, muchas veces, descontrolada.

Así pues, la descarga de toda comunicación puede realizarse de diferentes maneras, ya sea de forma síncrona o asíncrona, pudiendo agruparse en dos grandes categorías para conseguir esta información:

- **Poll.** Periódicamente se pide información al servidor de terceros.
- **Push.** El servidor manda información al cliente cuando hay información disponible para él.

Estos dos principales paradigmas son los más usados en la historia de Internet y de las aplicaciones. Durante mucho tiempo sólo se utilizaron las soluciones de polling, ya que no se había desarrollado una tecnología de push eficiente y eran más fáciles de implementar (preguntar cada cierto tiempo al servidor si había datos nuevos, manteniendo un pequeño intervalo en el cual cliente y servidor no estaban en sincronía). Sin embargo, ahora que la tecnología se ha ido desarrollando y se han creado nuevos modelos, el uso de notificaciones push está en auge, sobre todo gracias a los sistemas operativos móviles como Android o iOS, así como con nuevos estándares Web que permiten mantener sockets abiertos entre dos equipos diferentes para el envío de datos de forma bidireccional, como WebSockets¹ y Server-sent Events².

Sin embargo, el método de *polling* no es recomendable por varios motivos, debido a que hace que las redes se colapsen con su uso o bien que éste

¹<http://tools.ietf.org/html/draft-ietf-hybi-thewebsocketprotocol-17>

²<http://dev.w3.org/html5/eventsource/>

sea muy ineficiente. El primero de ellos es porque usa un elevado número de conexiones hacia el servidor de terceros para preguntar si hay datos. Además, esas conexiones no siempre son interesantes de realizar, puesto que puede que realmente no haya datos, pero al preguntar si los hay, se colapsa tanto el servidor como la red con paquetes innecesarios, que no tienen carga real alguna. Y es la combinación de estos dos factores (realizar conexiones periódicas, y que puede que no tengan nada) lo que indica que debería haber otra tecnología de recogida de información más eficaz y amigable con las redes.

Es por esto que los métodos de recogida de información basados en tecnologías push son ampliamente usados para la descarga de datos de forma asíncrona, ya que sólo piden los datos cuando realmente son necesarios. Sin embargo, aunque estos métodos a primera vista puedan verse como más eficientes que los de polling, también cuentan con varias desventajas y pese a ser realmente válidos para conexiones estables y sin limitaciones, son muy poco efectivos y muy perjudiciales para las redes móviles, como las de telefonía.

Resumen

La creación de un sistema de notificaciones push, con un API sencilla y fácil de utilizar, transparente para los desarrolladores de aplicaciones, con el objetivo de que funcione en diferentes tipos de redes de datos, como celulares o ADSL, en el dominio de paquetes o de circuitos, y haga un uso mínimo de recursos y de batería, utilizando todas las peculiaridades de cada servicio para que sea la mejor solución en cada situación.

Objetivos

- Explicar por qué es un problema usar cualquiera de las dos soluciones planteadas al inicio del capítulo en dispositivos móviles y redes celulares, viendo cuál es el problema inicial, las singularidades de dichas redes móviles y cómo pueden mejorarse.
- Creación de una plataforma de mensajería push que sea amigable para los dispositivos móviles con las consideraciones expuestas en el punto anterior: que no use demasiada batería, que sea interesante para las operadoras móviles y que reduzca el consumo de ancho de banda y de señalización en la red celular.
- Implementación de este mecanismo de push, no sólo en su parte servidor si no también en su parte cliente. Este desarrollo está englobado en la creación del sistema operativo para móviles Firefox OS.

- Estandarización al organismo W3C, para que lo implanten todos los navegadores y haya una mejor unión entre las redes de telefonía móvil y el mundo de Internet, para el intercambio de datos de forma eficiente.

Descripción del servicio

La plataforma del servidor de notificaciones se encargará de enviar mensajes push (pequeños mensajes, como conversaciones de chat, una estructura JSON sobre la descripción de un partido de fútbol...) a terminales que están dentro de las redes móviles, exponiendo unas APIs claras y sencillas, tanto para el desarrollador de la aplicación web como para el creador de los sistemas operativos.

El principal objetivo de crear este servicio es usar de una forma mucho más eficiente los recursos de las redes móviles disminuyendo el uso de la batería, a la vez que se reduciría el tráfico de señalización en la red de telefonía móvil y se permitiría a las empresas de telecomunicaciones ahorrar medios y ofrecer un mejor servicio a sus usuarios.

Capítulo 2. Estado del arte

En este capítulo se enumerarán las diferentes tecnologías que se han usado a lo largo del tiempo para mandar mensajes push a los dispositivos móviles. Algunas de ellas han sido creadas por los operadores, por lo que tienen un carácter de buen uso de las redes móviles. Sin embargo, ninguna de ellas se ha impuesto sobre las creadas por las empresas de software o de Internet por varias razones que se mostrarán más adelante.

Pero antes de comentar las diferentes propuestas que se han realizado en los últimos años, es necesario saber cómo funcionan las redes móviles a nivel radio y cómo los operadores tienen montada su infraestructura móvil.

Redes móviles: funcionamiento y problemas

Para crear un servicio que sea especialmente interesante para las redes móviles y que use pocos recursos, ancho de banda y en general, que cualquier operadora pueda pensar en él, hay que tener en cuenta los siguientes aspectos: el funcionamiento de las redes móviles que están desplegadas en el mundo, los detalles de las implementaciones, por qué se han hecho e incluso los acuerdos entre la capacidad de las líneas, los mensajes de señalización y otros muchos aspectos para que el servicio esté equilibrado entre los usuarios de las redes y los propietarios de estas.

Así pues, para conocer en profundidad el gasto de batería en los dispositivos móviles durante la conexión y por qué no se hace un buen uso de todos los recursos que aportan las redes móviles, es necesario empezar planteando un escenario satisfactorio que guste a todas las partes.

Redes LAN públicas o privadas

En primer lugar, hay dos tipos de redes principales para identificar los terminales dentro de las redes de telefonía, basadas en IP:

- **IPv4.** Son las direcciones IP más usadas en la actualidad y estandarizadas en el año 1981, en el RFC-791¹. Se caracterizan por tener 2^{32} direcciones disponibles para asignar, equivalentes a 4.294.967.296 equipos. Sin embargo, el principal problema es que son demasiado pocas para todos los diferentes dispositivos que hay conectados en la actualidad a Internet y requieren de técnicas (como el NATting) para poder acomodar más equipos, con problemas que se analizan más adelante.

¹<http://tools.ietf.org/html/rfc791>

- **IPv6.** Es la evolución directa de IPv4, estandarizada como RFC-2460², pero sin compatibilidad hacia atrás. De esta forma, los equipos que funcionan sólo con IPv6 no son capaces de llegar a dispositivos equipados sólo con IPv4 y para que puedan hablar entre ellos, hay que realizar túneles o conversiones entre ambos tipos de direcciones. Esto cuenta con una serie de ventajas sobre IPv4, la principal es que el número de dispositivos con IP distinta y no repetida aumenta de forma dramática (hasta el punto de llegar a 340 sextillones de direcciones, o 2^{128}) y que, en un futuro, relegará a IPv4 a un segundo plano.

Además de los tipos de direcciones, ya sean IPv4 o IPv6, existe el problema de si dichas IPs están en redes públicas o privadas, las cuales limitan el acceso que tienen hacia otros dispositivos dentro o fuera de dichas redes. Esto puede ser beneficioso en algunos casos y realmente perjudicial en otros.

Así pues, existen estos dos tipos de redes, que se explican mediante:

- **LAN pública.** Los equipos conectados a este tipo de redes poseen una dirección IP (ya sea de versión 4 ó 6) que les permite ser visibles para el resto de dispositivos conectados a Internet y, en determinada medida, accesibles sin restricciones. Es decir, cualquiera puede saber si un usuario está en Internet y mantiene una conexión directa con los equipos a los que se quiere conectar (servidores web, de correo, de mensajería instantánea), y ellos pueden contactar con dicho usuario aunque no tengan conexión, puesto que es visible.
- **LAN privadas.** Se encuentran dentro de un segmento de red que sólo es visible para otros equipos que lo compartan con nosotros. Esto hace que no se tenga una conexión directa con Internet y que haya que pasar por un dispositivo intermedio (normalmente un NAT) que permite salir a Internet no con nuestra dirección real, si no con la dirección que tenga dicho NAT y redirija nuestras peticiones hasta el exterior, a la vez que reconduce las respuestas desde el exterior hasta nuestro dispositivo.

Sin embargo, por diferentes motivos, ninguna de las soluciones es la panacea. En primer lugar, no contar con un control intermedio (normalmente un firewall) para protegerse de ataques, hace que tener una IP pública suponga una exposición de forma completa a cualquier dispositivo de Internet. Además, ser visible por todo el mundo, convierte en vulnerable cualquier petición que proceda del exterior, como escaneos de puertos, PINGs, o floods incontrolados. Es por esto, que la mayoría de operadores de telefonía móvil que disponen de servicios de datos, tratan de no dar IPs públicas a sus dispositivos para evitar estos problemas, que

²<http://tools.ietf.org/html/rfc2460>

al final llevan a un mayor uso de batería debido ya que reciben más datos, detalle que se explicará en la siguiente sección.

Pero a su vez, tener un direccionamiento privado tampoco ayuda en demasía a las redes móviles. No ayuda puesto que, aunque en algunos casos sea útil como en el momento de recibir notificaciones push, el inconveniente suelen ser problemas de seguridad o gasto de batería. Tener un direccionamiento privado nos convierte en visibles sólo para nuestros vecinos en nuestra red privada y dado que la IP interna es diferente a la externa, cualquier intento de conexión desde la red pública (Internet, por ejemplo) no llegaría a nuestra IP, ya que diferiría de la asignada para salir a Internet.

También se podría pensar que debido a que IPv4 tiene un rango de direcciones muy limitado, el uso de redes privadas con NATing está muy extendido, algo que es correcta. Pero también podría pensarse que usando IPv6 todos los problemas estarían resueltos y que al haber tantas posibilidades de direccionar dispositivos, el uso de mecanismos de push se irá perdiendo. Sin embargo, hay que tener en cuenta, como queda expresado anteriormente, que las operadoras suelen desplegar sistemas de firewall, por lo que aunque sus clientes dispongan de terminales con IPv6 pública, y por lo tanto, visibles desde Internet, se bloquearían algunas conexiones para así evitar la congestión en la red y un uso elevado de batería. De esta forma, habría que crear reglas específicas en los firewall para permitir conexiones entrantes hacia los dispositivos por parte de determinados lugares.



Nota: Redes fijas en el hogar

No se ha hablado de redes en el hogar puesto que la diferencia no es tan apreciable y permiten mantener equipos fuera de Internet, como televisiones o sistemas DLNA. Además, es transparente para los usuarios, ya que el peso lo lleva el router, que además de proporcionar acceso a Internet, actúa como NAT, facilitando sistemas 2 en 1.

Estados del dominio de circuitos

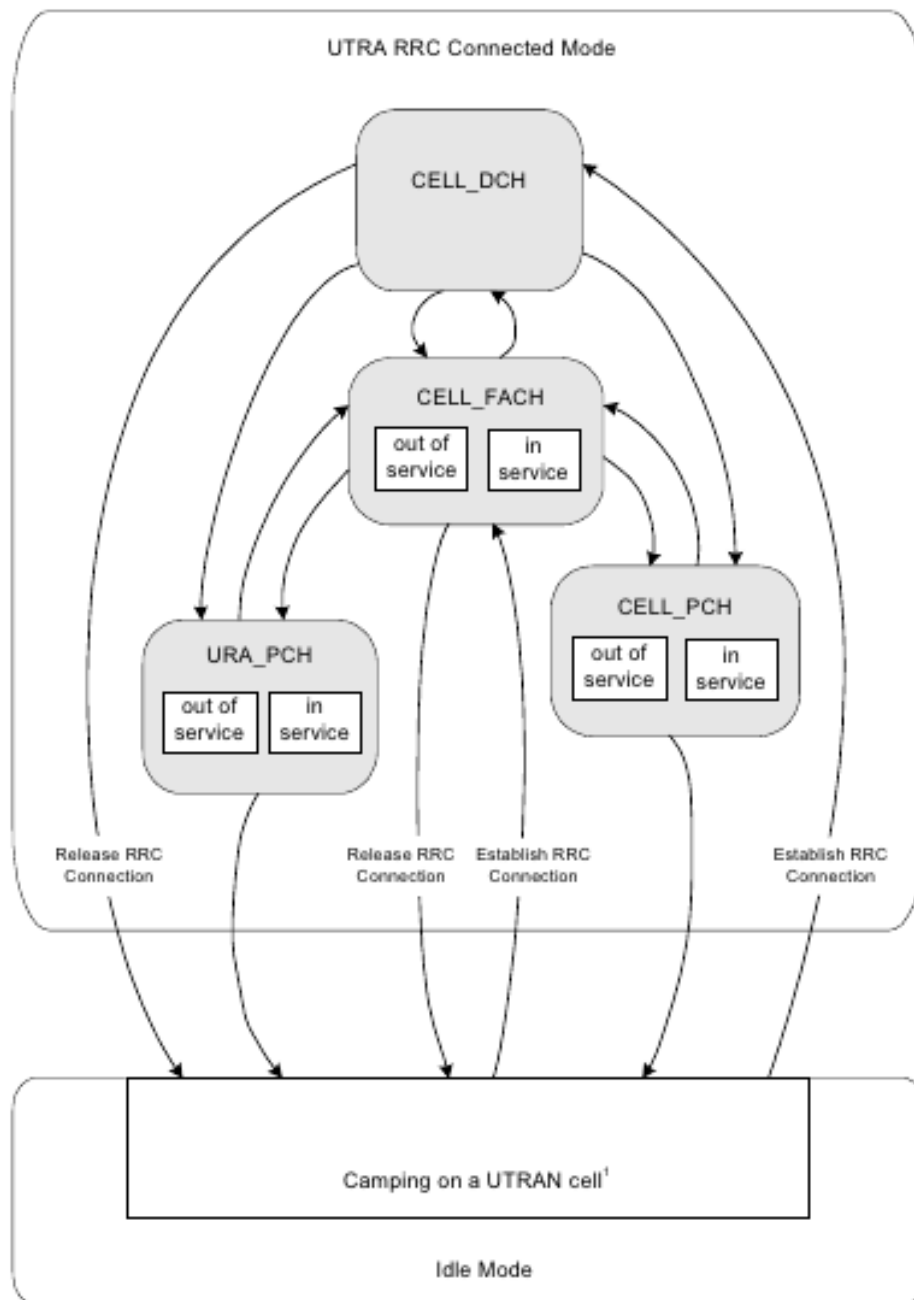
Como se ha comentado anteriormente, a continuación se procede a describir los estados radio más interesantes. En la especificación 3GPP TS 25.331³ se pueden ver todos los estados de la capa RRC (Radio Resource Control):

Para simplificar, sólo están listados los estados de la red 3G, actualmente la más usada.

³<http://www.3gpp.org/ftp/Specs/html-info/25331.htm>

- **Cell_DCH (Canal dedicado).** Cuando el teléfono se encuentra en este estado radio es porque está transmitiendo una gran cantidad de datos y la red le ha puesto en un canal dedicado. El tiempo de inactividad es muy corto para este estado, conocido como el temporizador T1, y suele variar entre 5 y 20 segundos. Esto quiere decir que si se alcanza este temporizador, el teléfono cambiará su estado radio a Cell_FACH.
- **Cell_FACH.** En este estado, el teléfono está conectado a la red móvil y usa un canal compartido con otros terminales. Suele transmitir un pequeño número de paquetes, que normalmente son keep-alives. El temporizador de inactividad en este estado es un poco mayor (del orden de los 30 segundos) y es conocido como el temporizador T2. Cuando se dispara, el dispositivo se mueve por la radio al estado Cell_PCH o URA_PCH (dependiendo del tipo de red).
- **Cell_PCH o URA_PCH (PCH: Paging CHannel) (URA: UTRAN Registration Area).** En este estado, el teléfono no puede mandar ningún dato, excepto información de señalización para poder localizar el dispositivo dentro de la red móvil. En este caso, la conexión RRC está establecida y abierta, pero apenas es usada. En dicho estado, el dispositivo informa a la red cada cierto tiempo de que el dispositivo cambia de una celda (o sector) a otra, para que la red sea capaz de conocer exactamente en qué BTS (Base Transceiver Station, o nodo B) se encuentra ofreciendo servicio al dispositivo. El temporizador T3 define el tiempo máximo que un terminal puede estar en un estado PCH. Este temporizador es mayor que T1 y T2 y depende de cada operador de red. Cuando se dispara, el dispositivo se mueve al estado IDLE, por lo que cualquier dato que quiera transmitir el dispositivo le llevará unos dos segundos para restablecer el canal y una gran cantidad de mensajes de señalización.
- **RRC_IDLE.** Es el estado más económico en el que puede estar la radio, ya que prácticamente está dormida. También, se ha de señalar que el modem está escuchando los datos de célula, por lo que cada vez que detecta que el usuario cambia de un LAC (Localization Area Code, un grupo de múltiples BTS o NodosB) a otro, el dispositivo cambiará el estado PCH para informar a la red de su nueva ubicación. Así pues, cuando el dispositivo se encuentra en este estado, puede pasar a un estado más activo de red. Esta también sabe en qué LAC se encuentra moviéndose el dispositivo, por lo que si la red necesita informar al dispositivo de que está ocurriendo algo (una llamada, un SMS, alguien quiere mandarle datos), se manda un mensaje de broadcast, llamado Paging por toda la red LAC de BTS para localizar el dispositivo.

El siguiente esquema representa los diferentes estados radio, ordenados por consumo energético en el dispositivo.

Figura 2.1. Esquema de estados radio

En definitiva, el estado más interesante para mantener un dispositivo móvil en la red celular es el modo Idle, puesto que no tiene una conexión de red abierta constantemente y el uso de recursos es de alrededor 100 veces menos entre este modo y el de máxima excitación (CELL_DCH), permitiendo que la batería dure más tiempo, haya menos señalización de red y se ocupen menos recursos, tanto en nuestro dispositivo como en la red.

Entonces, como se puede ver, se quiere mantener el dispositivo durante el máximo tiempo posible en este estado, por todas las razones expuestas anteriormente.

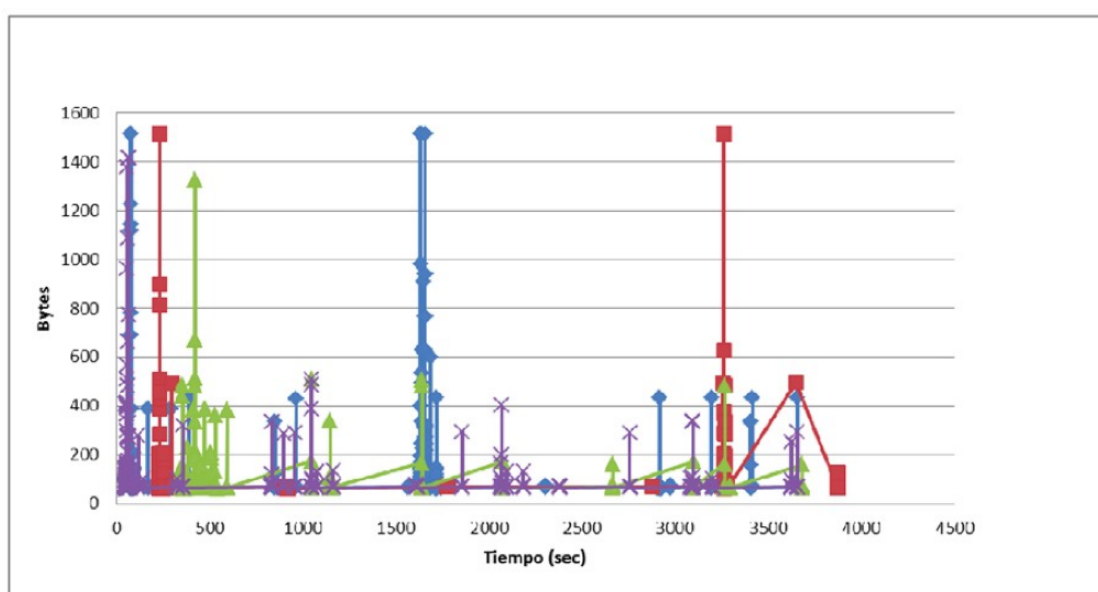
En las figuras siguientes, se muestran combinaciones de aplicaciones muy conocidas actualmente (como mensajería instantánea o VoIP) junto a las pruebas realizadas en diferentes terminales móviles, con los sistemas operativos más usados en estos momentos.



Nota

Por cuestiones legales, no pueden mostrarse los nombres de las aplicaciones y los terminales usados para generar estas gráficas. Pruebas realizadas por los laboratorios de Telefónica de España, S.A.

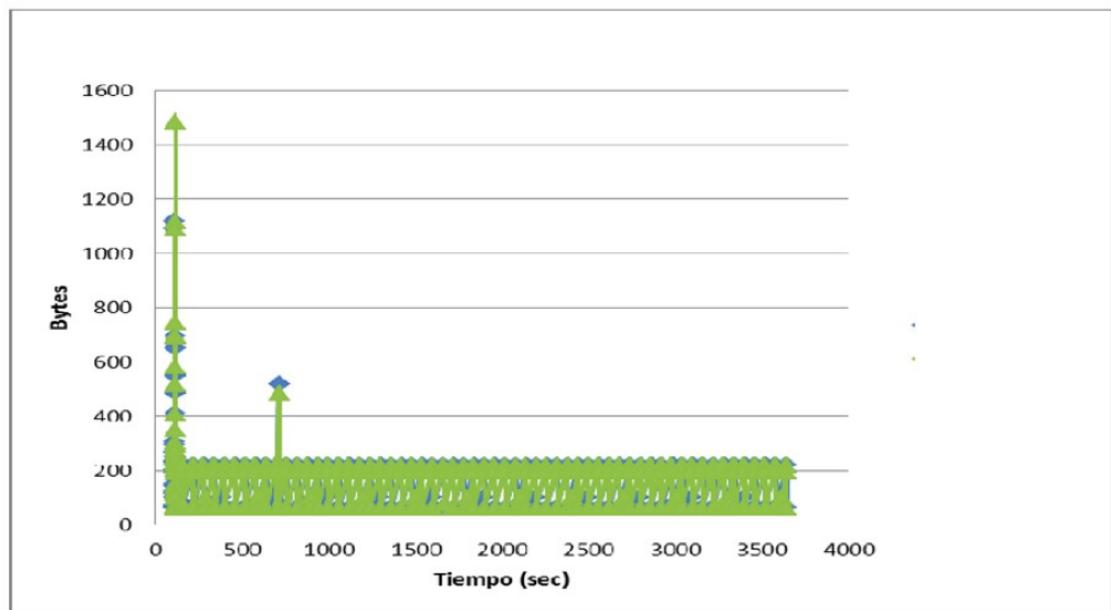
Figura 2.2. Aplicación con demasiado tráfico de control, 1



En la figura anterior, se puede ver el caso de una aplicación que utiliza muy mal los recursos de la red. Abre una conexión que mantiene de forma permanente con los servidores de su aplicación, ya sea para indicar actividad o presencia, recepción de mensajes de control u otras funciones.

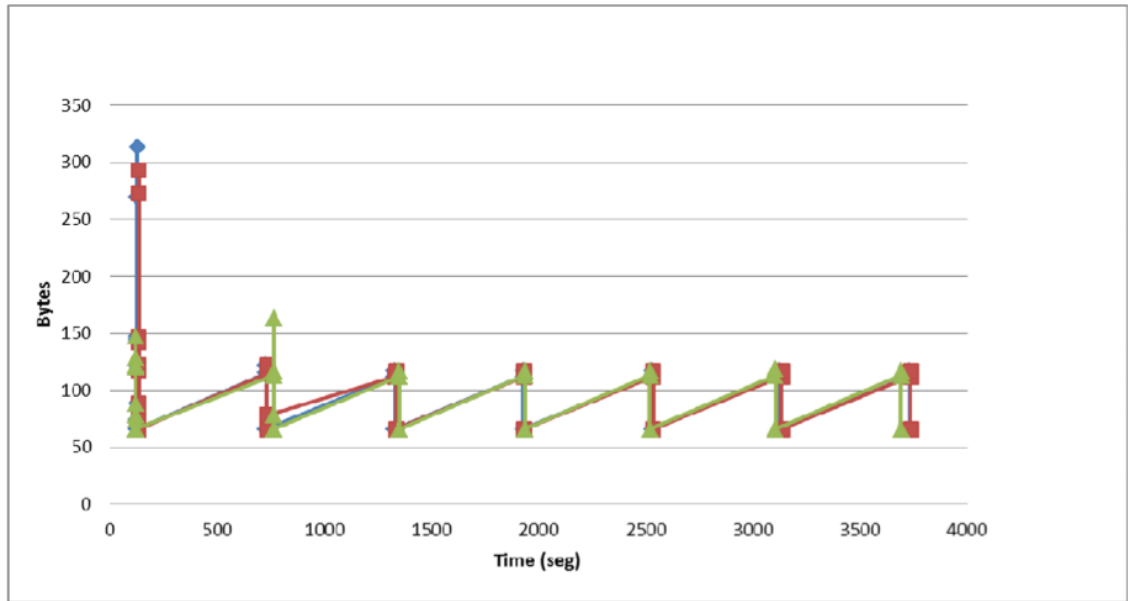
Cada color es un terminal diferente, que pueden tener diferentes sistemas operativos y versiones, y se puede observar el comportamiento tan errático que tiene a lo largo del tiempo. Manda mensajes en tiempos totalmente aleatorios, pero próximos en el tiempo, incluso del orden de kilobytes.

El problema principal es que el modem radio no puede quedarse en un estado de bajo consumo energético, debido a estos cambios totalmente aleatorios. Cada vez que una aplicación quiere enviar datos, tiene que restaurar la conexión de radio móvil, actualizar hasta el estado necesario para poder empezar a transmitir, y empezar a enviar y recibir datos.

Figura 2.3. Aplicación con demasiado tráfico de control, 2

Si la primera gráfica indicaba un comportamiento errático de una aplicación, con unos envíos de datos totalmente aleatorio, esta indica todo lo contrario, con un comportamiento muy estable, pero tremendamente perjudicial para la red.

Está todo el rato enviando datos de control, del orden de unos 256 bytes, y durante intervalos muy cortos de tiempo, del orden de segundos. Esto lo que provoca es que el modem radio del terminal sea incapaz de llegar a los temporizadores mostrados anteriormente para bajar su nivel de actividad, y por tanto, su nivel de uso de batería. Es probablemente mucho más perjudicial para las redes y las baterías de los móviles que la aplicación anterior.

Figura 2.4. Aplicación con poco tráfico de control

En esta gráfica, se puede ver una aplicación que hace un uso más coherente de los recursos de las redes móviles, haciendo `PINGS` cada 10 minutos (o 600 segundos) para informar de presencia o cualquier otro dato de control.

Sin embargo, también tiene un problema, y es el uso de `PINGS` para informar de diferentes hechos en la aplicación. Podría ser aceptable tener una sola aplicación en el móvil que haga esto, pero sin embargo, los usuarios suelen tener varias instaladas: alguna aplicación de mensajería instantánea, otras de VoIP, clientes de lectura de RSS, información sobre correos nuevos... El problema viene cuando cada una realiza exactamente esto, que al final se está haciendo `PINGS` cada minuto o menos, por lo que la suma de todas las aplicaciones hace que cualquier precaución tomada de forma independiente por cada una de ellas, se evapore.

Operadores

Una vez analizados los problemas de la red móvil, cómo se comporta y cuál es el mejor estado de radio para mantener a un terminal para que el uso de batería y de recursos de red sea el menor posible, se procede a explicar las soluciones actuales de push, tanto las propuestas por las propias operadoras de telefonía, ya sea por parte del GSMA⁴ como por OMA⁵, dos de los organismos internacionales en la estandarización de protocolos para telefonía móvil.

⁴<http://www.gsma.com/>

⁵Open Mobile Alliance: <http://openmobilealliance.org/>

WAP Push

Históricamente, los operadores móviles han ofrecido y ofrecen mecanismos reales para notificaciones push, conocidos como WAP Push. Este mecanismo es muy eficiente con la red, ya que incluso permite que no se tenga una conexión de datos abierta (con su contexto PDP), dado que no funciona bajo el dominio de paquetes, si no de circuitos. Su funcionamiento es muy similar a la recepción de un SMS o de una llamada de teléfono, en el cual el operador busca si el móvil está conectado, en qué área se encuentra disponible y posteriormente hace un broadcast a la célula para despertar al teléfono (mediante un mensaje de PAGING) y decir que tiene algo disponible para él.

Lo interesante de este sistema es que es algo integrado de forma perfecta en la red, puesto que no son más que mensajes de SMS especialmente compuestos, por lo que usa los "restos" del espectro móvil para ser transmitidos y realmente no conllevan un gran coste para las operadoras.

Sin embargo, tiene un problema fundamental: las operadoras quisieron cobrarlo y aún hoy lo cobran. Para ellas, apenas tiene un coste económico, sin embargo, supone un gran ahorro de batería y de recursos de red si se hubiera ofrecido desde un primer momento de forma gratuita y no se hubiera obligado a crear soluciones propietarias que dañan las comunicaciones, como se hablará posteriormente en la sección de Internet.

Su funcionamiento es muy sencillo. El servidor de terceros quiere enviar cualquier tipo de mensaje binario (aplicaciones, imágenes, vídeos, simple texto), por lo que en primer lugar se sube a un servidor que está en Internet, o bien en la red privada del operador, y después se manda un SMS con una estructura XML determinada que contiene una URL con la dirección que hay que descargar, que se compila y se envía como un SMS binario. Posteriormente, el mensaje es interpretado por el sistema operativo del móvil en cuestión y abre el navegador apuntando a la dirección en concreto que se había señalado en el mensaje WAP Push, iniciando la descarga.

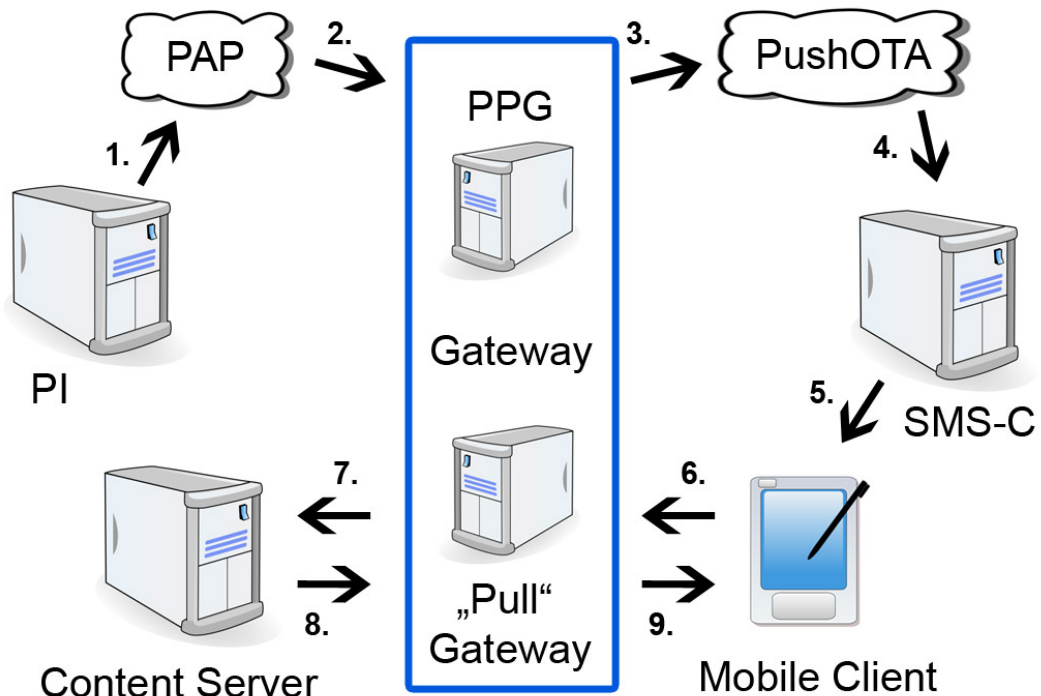
Como se puede observar, esto es un problema, puesto que está pensado para la descarga de datos desde Internet y obliga a que se abra el navegador de forma predeterminada para la descarga del contenido.

Visto el aspecto negativo, cabe destacar los positivos. El primer de ellos es que no utiliza un canal de datos dedicado en el dominio de paquetes, si no en el dominio de circuitos, y es una tecnología nativa en las redes móviles. Además, al ya haber despertado al teléfono, y por lo tanto la radio, para descargar el SMS, el estado de consumo es alto, por lo que una segunda conexión a un servidor en Internet para descargar los datos

reales no es tan perjudicial, sabiendo que sólo se conectará cuando haya datos que recoger, y no cada cierto tiempo.

Otro gran punto positivo es que pueden coexistir tanto el WAP Push como otra plataforma de push, incluso aunque no funcione bajo la red de telefonía, como el Wi-Fi. Por ejemplo, se podría recibir una notificación push por el método WAP Push y recoger los datos vía Wi-Fi si está disponible, por lo que normalmente la descarga será más rápida, utilizará menos energía y bajará aún más la señalización necesaria para poder descargar, puesto que utiliza la red local y no la red proporcionada por el operador.

Figura 2.5. Funcionamiento WAP Push



Se puede explicar el funcionamiento de WAP push:

1. PI (Push Initiator), que es normalmente una aplicación que corre en un servidor web, se comunica usando el protocolo PAP (Push Access Protocol), que no es más que un XML usado para expresar las instrucciones de entrega del mensaje.
2. Este protocolo hace que el PI se entienda con el PPG (Push Proxy Gateway), que es el encargado de recibir el mensaje inicial. Este PPG es el principal actor para WAP Push, del que se han recogido varias de sus ideas para el servidor de push que se desarrolla en este proyecto. Aparte de recibir el mensaje, es necesario que lo transforme o adapte

para el receptor en concreto, guarde el mensaje por si el usuario no está conectado e incluso una conversión entre la dirección de recepción y el destinatario final.

3. Una vez procesado el mensaje, el PPG se comunica mediante PushOTA (Push On-The-Air) con el servidor central de SMS-C para pasar el mensaje (o mensajes) SMS que compondrán el contenido final. El protocolo PushOTA puede ser de varios tipos: una interfaz OTA-HTTP o una conexión OTA-WSP.
4. El mensaje es recibido por el SMS-C, que mirará a qué número de teléfono va dirigido, se encarga de enrutarlo, etcétera.
5. El mensaje se envía al dispositivo móvil del cliente usando la especificación normal de entrega de mensajes cortos SMS. Después, el sistema operativo del móvil interpreta el mensaje, viendo que es un mensaje WAP push y abre la página web que contiene el paquete binario que el proveedor quiere entregar. Esta apertura del navegador puede ser automática o bien aceptada por el usuario.
6. El navegador carga la dirección (que puede pasar por el PPG, pero no es obligatorio).
7. (el mismo paso que el anterior, en el caso de que exista proxy).
8. El servidor de terceros donde realmente está alojado el contenido devuelve la información requerida por la dirección que ha abierto el navegador del cliente. Este es el paso final y el momento en que el usuario cuenta con los datos requeridos en un primer momento.
9. (pasando por el proxy intermedio si hiciera falta).

De este modo, se puede observar que el mecanismo de WAP push está diseñado para funcionar muy en sintonía con la propia red móvil. Esto supone que la transmisión sea muy sencilla y con un bajo coste para las operadoras, utilizando las redes de forma nativa para transmitir los datos necesarios de cara a llegar al contenido final. Además, contribuye a que los dispositivos puedan estar en modo de espera, con poco uso de batería para ser despertados en el caso de que tengan notificaciones, de la misma manera que se realizan los SMS y las llamadas.

Internet

Una vez vistas las diferentes plataformas push creadas por los organismos de los operadores de telefonía móvil, el siguiente paso es explicar otros mecanismos, caminando hacia tecnologías más modernas

creadas por empresas radicadas principalmente en Internet. Estas tecnologías son relativamente guardan importantes similitudes entre ellas, pero difieren en las cuotas, autenticación... Sin embargo, la idea de cómo mantener el canal de comunicación entre los dispositivos y el servidor de notificaciones es muy semejante.

A continuación se procede a una breve explicación de cada una de estas tecnologías.

GCM: Google Cloud Messaging

GCM o Google Cloud Messaging⁶ es un sistema de notificaciones push creado por la empresa Google para su sistema operativo Android. Su finalidad es la misma que la del resto de sistemas: entregar notificaciones o mensajes que un servicio de terceros (o incluso la misma plataforma de Google) a los usuarios empleando un dispositivo de una manera fácil, ordenada y controlada.

Lanzado en julio de 2012, es uno de los sistemas más grandes, puesto que está presente en todos los dispositivos Android que dispongan de una cuenta de Google asociada, prácticamente la mayoría. El número de usuarios de este servicio no es especificado por Google, pero podría ser similar o inferior al número de dispositivos Android en el mercado, que supera los 500 millones⁷.

Google Cloud Messaging for Android (GCM) es un servicio que faculta enviar datos desde su servidor a los usuarios con dispositivos Android. Puede tratarse simplemente de un pequeño mensaje que diga a la aplicación que hay nuevos datos para ser descargados desde el servidor (por ejemplo, que un amigo ha subido un nuevo vídeo) u otro mensaje que contenga hasta 4KiB de información (por lo que las aplicaciones como las de mensajería instantánea pueden consumirlo directamente).

El servicio GCM maneja todos los aspectos del encolamiento de los mensajes y entrega a la aplicación Android determinada que se puede ejecutar en el dispositivo.

—GCM: Google Cloud Messaging for Android

⁶<http://developer.android.com/google/gcm/index.html>

⁷A fecha de 12 de septiembre de 2012, http://news.cnet.com/8301-1035_3-57510994-94/google-500-million-android-devices-activated/

Figura 2.6. GCM logo



Características. GCM tiene como características más interesantes:

- La aplicación a la que tiene que llegar la notificación no tiene por qué estar abierta. El servicio GCM se encarga de despertarla y de que maneje la notificación.
- No hay ninguna interfaz para administrar la notificación. GCM simplemente pasa datos en crudo desde que la notificación es recibida hasta que es entregada. La aplicación final es la goza de la lógica necesaria para saber qué hacer con la notificación en todo momento. Por ejemplo, puede que la notificación sea "hay nuevos datos para descargar", por lo que no muestre ninguna interfaz, o bien un nuevo mensaje de un amigo, que sí es importante mostrar.
- Requiere que la versión de Android instalada sea la 2.2 y además es necesario que posea la aplicación de Google Play (el antiguo Market) configurada correctamente, es decir, una cuenta de Google asociada. Pero no obliga a que la aplicación que reciba la notificación esté instalada por Google Play, si no que puede tener cualquier origen.
- Usa una conexión permanente con los servidores de Google para recibir las notificaciones y mensajes de control.
- Requiere que los servidores de terceros que envíen notificaciones estén registrados en la plataforma. Esto significa que Google puede revocar en cualquier momento el envío de notificaciones por parte de un desarrollador o empresa, simplemente no permitiendo el identificador único que se les proporciona en un primer momento.
- La identificación de la aplicación es de forma única. Esto significa que para cada aplicación que quiera recibir notificaciones push, en cada dispositivo, tiene un número único que la identifica de forma única dentro del servicio. Esta identificación única la proporciona el servidor de GCM a la aplicación, que a su vez tiene que enviarla al servidor de terceros para que se destine concretamente a esa instancia.

El diagrama principal del envío y entrega de un mensaje de una forma muy abstracta es:



Nota

Se supone que la aplicación instalada tiene un ID de registro que permite recibir notificaciones, además de que el servidor de terceros ha guardado ese ID y el servidor tiene una clave de API que le permite identificarse como emisor de notificaciones.

1. El servidor de aplicación envía un mensaje a los servidores de GCM.
2. El servicio GCM encola el mensaje y lo guarda en el caso de que el dispositivo al que se tiene que entregar esté desconectado.
3. Cuando el dispositivo se encuentre online (puede que ya lo esté, o puede que no), el mensaje es mandado por un canal al dispositivo, que lo parsea el sistema operativo (Android, en este caso), y despierta a la aplicación destino en el caso de que esté cerrada o la entrega directamente.
4. La aplicación procesa el mensaje y realiza cualquier tipo de evento relacionado con ella: sincronizar datos de fondo, mostrar una alerta...

GCM es un sistema que no es obligatorio para las aplicaciones Android publicadas en el Market de Google. Sin embargo, su diseño le convierte en algo bastante eficaz para las redes móviles y le confiere una serie de ventajas que hacen que los servicios de terceros tengan que disponer de menos lógica y preocuparse en menor medida de la suerte que van a correr sus mensajes.



Nota

Cabe señalar, que antes de la introducción de GCM para los desarrolladores de Android (en julio de 2012), había un sistema muy similar para teléfonos Android, llamado Cloud to Device Messaging, o más conocido como C2DM⁸

Ventajas. GCM tiene una serie de ventajas sobre su predecesor CD2M (no explicado puesto que ha sido sustituido por esta versión) de las cuales habría que distinguir:

- Atesora un mejor uso de batería, ya que puede encolar mensajes para ser entregados cuando el dispositivo está con una conexión activa (y por lo tanto tiene la radio móvil activada). Además, no tiene por qué notificar al dispositivo siempre que le llegue una nueva notificación, si no que puede esperar a que pase a estar activo (por ejemplo, que se encienda la pantalla) para mandar los datos.

⁸<https://developers.google.com/android/c2dm/>

- El uso de datos y la transferencia es más eficiente, según las explicaciones del punto anterior. Esto lleva a un menor uso de batería y a que el usuario sólo pueda recibir notificaciones cuando realmente sean interesantes: cuando se usa el dispositivo y mantener en cola las no prioritarias.
- El API es más sencilla, puesto que hay menos pasos para poder usar el servicio y además, el código del cliente en las aplicaciones es más claro y fácil de implementar, sobrescribiendo algunos métodos de las clases de las cuales se extiende.
- La migración del servicio es muy simple y sólo hay que cambiar la URL de push para que apunte hacia otro servidor.
- Se eliminan las cuotas que había en el sistema anterior, permitiendo a los desarrolladores crecer de forma sencilla y sin nuevas peticiones. Igualmente, se puede consultar el estado de los mensajes en el perfil de Google.
- Permiten un mensaje (`payload`) de 4KiB, por lo que muchos de los datos que puedan ir por esta mensajería push no requerirán que las aplicaciones se conecten al servidor de terceros para recoger la información y realizar una segunda conexión.
- Posibilitan los mensajes multicast, es decir, se puede realizar un sólo envío de mensaje push, pero indicando todos los receptores en el momento de enviarlo para que se entregue a todos. Incluso múltiples emisores pueden enviar mensajes a una misma aplicación (por ejemplo, una red social que habla sobre un jugador de fútbol y a la vez información sobre el equipo en el que juega).
- Hay un tiempo de vida máximo de 4 semanas para cada mensaje. Si un mensaje viene sin tiempo de vida, se presupone que el tiempo es máximo, sin embargo, un mensaje muy interesante es del de un TTL (`time_to_live`) de 0 segundos, el cual se envía a los receptores sin guardarlo, de forma que si no está conectado el receptor (o los receptores), se pierde.
- Existe la posibilidad de sobrescribir mensajes antiguos (con la misma `collapse_key`), enviar mensajes sólo cuando el usuario esté activo (con `delay_while_idle`) e incluso comprobar si los mensajes van a ser entregados o no, dependiendo de la respuesta que aporte el servicio.

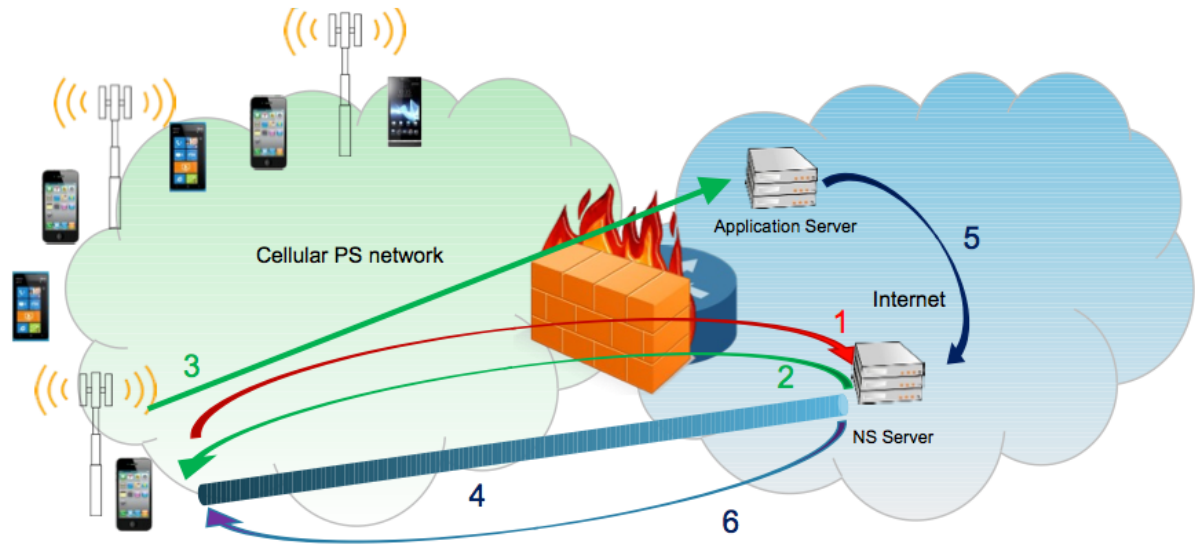
A pesar de la enorme capacidad de la que dispone este sistema, con todas las ventajas explicadas anteriormente, sigue teniendo pequeños fallos que hacen que no sea ampliamente usado, por ejemplo:

1. **No es obligatorio.** Por lo que los diferentes desarrolladores no están obligados a usarlo en sus aplicaciones y pueden crear otros sistemas

de push propietarios (como los que hacen la mayoría de aplicaciones de mensajería instantánea), lo que lleva a que la validez del sistema sea menor.

2. **Mantiene una conexión abierta.** Como se ha comentado, uno de los grandes problemas, junto a la señalización, es la de mantener canales abiertos sin hacer nada a la espera de que alguna vez atravesase un dato, en este caso, notificaciones.

Figura 2.7. Funcionamiento de GCM



Arquitectura. Así pues, se puede observar el flujo completo de trabajo de GCM:

1. En la primera utilización del servicio, el teléfono dotado con Android (la única plataforma que soporta GCM por el momento) tiene que disponer de una cuenta de correo válida en el sistema (de Gmail o de Google Apps), a la cual se asociará un identificador único de dispositivo. Asimismo, también se enviará una petición para recibir una cadena única que identificará a una aplicación en un dispositivo en concreto.
2. La aplicación que pidió la identificación, la recibe del servidor GCM.
3. La aplicación manda, por su propio canal de comunicación que haya decidido, dicho token único a un servidor de terceros, con algún tipo de identificación de usuario o de instancia que permitirá enviar contenido personalizado. Por ejemplo, si una persona registrada desea recibir notificaciones sólo sobre su equipo favorito, se tiene que enviar el dato que identifique a esa persona/usuario con el identificador enviado.
4. El dispositivo abre y mantiene una conexión TCP persistente con el servidor de GCM, mandando mensajes de keep-alives (que viene a ser un "sigo vivo") en intervalos de tiempo determinados. Estos intervalos

son, por defecto, de 28 minutos, pero varía dependiendo del país y de la operadora, e incluso de forma adaptativa.

5. En cualquier momento, el servidor de terceros (llamado Application Server) manda un mensaje push al servidor GCM con los parámetros que se especifican en el servicio, incluyendo el identificador de dispositivo recibido en el paso 2 y enviado en el paso 3.
6. GCM envía, en la medida de sus posibilidades, el mensaje al dispositivo a través del canal TCP abierto, ya sea de forma inmediata, o bien con los parámetros indicados en el paso anterior (como `delay_while_idle` o `time_to_live`).

Problemas. Sin embargo, a pesar de todas las nuevas características y todo lo que ha trabajado Google en este sistema, existen varios problemas que podrían llegar a ser bastante graves derivando en una no confianza hacia este sistema.

- Se necesita mantener siempre una conexión abierta, por lo que sigue ocupando espacio en la red móvil, y en concreto en el GGSN (GPRS Support Node), que es un hardware con un precio muy elevado donde se mantiene abierta cada conexión, hasta un límite de 65535.
- Esa única conexión requiere de keep-alives para mantenerse activa, por lo que cada cierto tiempo se estará cambiando de estado radio, con su correspondiente señalización, indicando simplemente que se está conectado.
- Es un punto único de fallo (SPOF: Single Point Of Failure) son los servidores de Google, por lo que si caen, la infraestructura también lo hace, y no hay posibilidad de cambiar a otra, dado que el protocolo es cerrado y propietario.
- Todos los mensajes van por los servidores de Google, por lo que puede constituir un posible riesgo de seguridad, más aún viendo las diferentes leyes para regular Internet que han salido desde Estados Unidos, donde Google tiene su sede.
- No hay garantía sobre el orden de entrega de los mensajes, por lo que algunos de estos enviados posteriormente podrían llegar antes que uno enviado previamente. De esta forma, no se pueden crear dependencias entre mensajes.

APNS: Apple Push Notification Service

APNS o Apple Push Notification Service es otro de los sistemas creado por empresas fabricantes de sistemas operativos para solventar el problema de notificaciones push, en este caso para dispositivos iOS, como el iPhone o el iPad.

Este sistema es muy similar en comportamiento a GCM pero con alguna variación, sobre todo en los tipos de mensajes que puede enviar los desarrolladores y en el tratamiento de las notificaciones por el usuario, que se expondrá más adelante.

Historia. Fue lanzado con la versión 3.0 de iOS, el 17 de junio de 2009, que correspondía . Ha sido mejorado con el paso del tiempo e introducido en los equipos de escritorio con la versión 10.7 de MacOS usando el centro de notificaciones, también presente en iOS 5 y posteriores. Su primera versión fue retrasada algunos meses debido a los problemas de escalabilidad que se estaban teniendo en el sistema por la gran avalancha de peticiones y uso que estaban dando los desarrolladores a este sistema⁹.

Características. Este sistema posee una serie de características que no tienen los demás métodos de notificaciones push, las cuales son:

- Usa un protocolo binario para el envío de datos, bajo una conexión TCP cruda. Esto hace que el tamaño de las notificaciones y de los mensajes de control sea más pequeño, pero a la vez sea más difícil de descifrar, puesto que el protocolo es cerrado.
- Su uso es muy amplio y obligatorio para todas las aplicaciones certificadas por Apple (todas aquellas que estén en su tienda de aplicaciones) en vez de implementar servicios de polling.
- Hay un canal de escucha de feedback, por lo que se puede saber si las notificaciones se han enviado correctamente o no. Además, también informa sobre de-registros de aplicaciones y dispositivos.
- Está controlado absolutamente por Apple, pudiendo censurar y eliminar dispositivos y servidores que estén haciendo un uso indebido de su infraestructura.

Ventajas. En primer lugar, se van a enumerar las diferentes ventajas que tiene sobre sus competidores:

- Mayor seguridad: las notificaciones recibidas para aplicaciones que están en segundo plano o cerradas, no son enviadas directamente, si no que se muestran en la barra de estado, y hasta que el usuario no pulsa y las acepta, no se envían a la aplicación, añadiendo un nivel más de seguridad, la aceptación del usuario.
- Funcional tanto en los dispositivos móviles, como iPad o iPhone, como en los sistemas de escritorio, tipo MacOS X, a partir de su versión 10.7, siendo el API el mismo para ambos sistemas.

⁹<http://www.engadget.com/2008/06/09/iphone-push-notification-service-for-devs-announced/>

- Permite tres tipos diferentes de notificaciones, que se adecuan a diferentes opciones visuales y sonoras:
 - **Badges.** Muestra una imagen o número en el icono de la aplicación.
 - **Alerta/Banners.** Muestra un aviso en el centro de la pantalla.
 - **Sonido.** Reproduce un sonido determinado.

Problemas. Sin embargo, hay una serie de desventajas que merece la pena reseñar:

- Protocolo muy cerrado y sin posibilidad de implementaciones de terceros, en ningún sistema operativo.
- El payload que pueden llevar las notificaciones es como mucho de 256 KiB, lo que hace que muchas aplicaciones necesiten una segunda conexión para recoger los datos reales, ya que pueden ocupar más de lo permitido.
- El servidor que envía las notificaciones debe tener abierta otra conexión con un servidor de feedback de Apple para recibir el estado de las entregas.
- Demasiada burocracia: creación de certificados ad-hoc, certificados SSL, es necesario registro, para probar se tienen que utilizar un dispositivo en concreto registrado con Apple...
- Sistema centralizado y propenso a que si hay un fallo, ningún dispositivo de Apple pueda recibir notificaciones.
- Cuotas muy estrictas: un mensaje por aplicación y hasta un máximo de 100 por dispositivo.

Como se puede observar, es un sistema bastante más limitado que GCM, con menores opciones, mayor control sobre lo que pueden hacer las aplicaciones y los desarrolladores, pero con la ventaja de que funciona en todos los sistemas operativos de Apple.

Thialfi: el futuro de Google

Thialfi es un nuevo protocolo diseñado y publicado por Google¹⁰ que trata de arreglar algunos de los problemas de los sistemas de notificaciones: el paso de mensajes completos a través de la red y el consiguiente inconveniente de infraestructura y escalabilidad que podría resultar de tener millones de terminales conectados a la vez mandando muchos mensajes a través de la red y en espera para ser entregados.

¹⁰<http://research.google.com/pubs/pub37474.html>



Nota

Sin embargo, cabe observar que este sistema no está funcionando en los terminales móviles que llevan Android, sino que se usa el ya explicado GCM, sólo con la intención de explorar nuevas posibilidades, así como la sincronización de contactos de manera experimental. Lo implementa Google Chrome en su versión de escritorio y móvil.

El principal cambio entre GCM y Thialfi es la vuelta de mensajes que no llevan contenido útil para la aplicación, si no que se basan en mandar números de versión sobre temas en concreto ("topics") a los que el usuario y las aplicaciones se suscriben para mantenerse en sincronía.

Características. Por consiguiente, Thialfi guarda una serie de mejoras con respecto a los otros sistemas:

- **Escalable.** Puede tener constancia de millones de clientes y de objetos, ya que el tamaño para ambos es pequeño, sin necesidad de guardar datos finales, si no simplemente números de versión.
- **Rápido.** Según las pruebas y el paper original, la entrega de los mensajes, en este caso mensajes con números de versión, se realiza en menos de 1 segundo desde que el servidor de terceros manda un mensaje hasta que el cliente lo recibe.
- **Seguro.** Incluso si un centro de datos se cae, el sistema es capaz de recuperarse, volviendo a registrar los diferentes dispositivos y eventualmente regresando a un estado correcto.
- **Fácil de usar.** Las APIs son sencillas y han sido implementadas en diferentes servicios utilizados por Google, como la sincronización de datos de Chrome, contactos y en Google+.

Entonces, tal y como se ha expuesto en los puntos anteriores, todo son objetos que poseen identificadores (ID) únicos y números de versión que incrementan de forma monótona en cada actualización.

Ventajas.

- **Menos datos guardados.** Es una de las ventajas, pero a la vez es un gran inconveniente. Sólo hay que guardar los objetos que contienen un identificador y un número de versión, por lo que las bases de datos (BigTable, que usan en Google) necesitan menos espacio para poder funcionar, al no guardar mensajes (en GCM el `payload` máximo es de 4KiB). Además, sólo es necesario guardar el último mensaje para cada aplicación-objeto (básicamente es aumentar el número de versión),

mientras que en un sistema como GCM se requiere guardar múltiples mensajes.

- **No hay fugas de datos privados.** Dado que las notificaciones no tienen mensaje, no hay riesgo de fuga de datos privados, puesto que lo único que se transmite son números de versión, que no dicen nada sobre el contenido.
- **Entrega siempre actualizada.** La entrega de los mensajes está asegurada, no porque siempre se tenga la última versión en un momento determinado, si no porque seguro que se va a conocer esa versión en un momento futuro. Esto quiere decir que lo único que reciben los clientes es un número de versión de un determinado objeto, que siempre es el máximo posible. Así pues, si un cliente se ha quedado desactualizado en, por ejemplo, la versión 17, y ha llegado la 18 y no ha podido ser entregada, cuando llegue la versión 19 el cliente sabrá que tiene que pedir tanto la 18 como la 19, para así estar en sincronía con el servidor, de manera que en cualquier momento lo único que almacena un servidor o un cliente es un determinado objeto, con un ID en una versión v determinada.
- **Buena recuperación ante errores.** Thialfi tiene una serie de mensajes de control que permiten recuperar el estado de los registros en caso de fallo. Por ejemplo, si el servidor pierde todos los registros, se envía un mensaje a los clientes para que ejecuten `ReissueRegistrations()`, y se inicia un protocolo para devolver los registros.
- **Fácil implementación del cliente.** El número de líneas requeridas para la implementación en el cliente es muy bajo, debido a su uso basado en eventos y su fácil lógica, en la siguiente figura se pueden ver algunos ejemplos:

Figura 2.8. Líneas de código para implementar Thialfi

Application	Language	Network Channel	Client Lines of Code (Semi-colons)
Chrome Sync	C++	XMPP	535
Contacts	JavaScript	Hanging GET	40
Google+	JavaScript	Hanging GET	80
Android Application	Java	C2DM + Standard GET	300
Google BlackBerry	Java	RPC	340

Como se puede observar, varía dependiendo del lenguaje de programación, desde las 40 u 80 en Google Contacts y Google+ hasta las 535 bajo un protocolo XMPP en Chrome Sync.

Además, lo interesante que se ve en esta figura anterior es que Thialfi no obliga un único protocolo de aplicación, si no que puede funcionar sobre varios, como conexiones HTTP, XMPP o incluso RPC.

Problemas. Sin embargo, aunque aporta muchas ventajas, como las que están enumeradas anteriormente, cuenta con problemas de cara a ser un buen servidor para datos en la red móvil

- **Doble conexión para recuperar datos.** En un sistema en el que los servicios de terceros pueden enviar datos de payload, obliga a que el servidor de notificaciones sea más complejo, por el hecho de guardar datos. Sin embargo, esto puede significar que los clientes no tengan que abrir una segunda conexión para recuperar los datos necesarios debido a que el propio mensaje de notificación ya lo tiene. Esto, a pesar de parecer un problema menor, constituye uno de los grandes inconvenientes en las redes móviles, donde realizar una conexión es una operación muy costosa (en términos de tiempo y de señalización) y la descarga de los datos es lenta.
- **Mayor complejidad en el servidor.** Esto se debe a que requiere guardar los cambios o diferencias para que cuando se conecten los diferentes clientes, puedan recuperar los datos que pidan, como por ejemplo, los datos de la versión 17. Sin embargo, puede haber casos en los que los clientes no vayan a pedir la última versión, si no la versión 16 y 17. Por lo que el servidor debe mantener esas diferencias de algún modo (en formato delta o con mensajes completos).
- **Mayor complejidad en las aplicaciones.** Esto es debido a que las aplicaciones no van a procesar el mensaje directamente, si no que tienen que discernir cuál es el número de versión, hacer una llamada al servidor de terceros y recibir los datos. Además, se ha de tener en cuenta que pueden solicitarse diferentes versiones para acabar en sincronía con el servidor, lo que aumenta la lógica para ver qué versión está dando el servidor y cómo se puede manejar con versiones más nuevas.
- **Servidores de aplicación más potentes.** El primer problema es la necesidad de que todos los clientes, una vez recibido un número de versión, tengan que conectarse obligatoriamente con el servidor de la aplicación para descargar los datos, lo que podría conllevar muchísima capacidad de cálculo (si por ejemplo los datos hay que calcularlos al vuelo y dependen de cada petición) y de ancho de banda, así como

de medidas para evitar que muchos usuarios se conecten a la vez y puedan colapsar los sistemas.

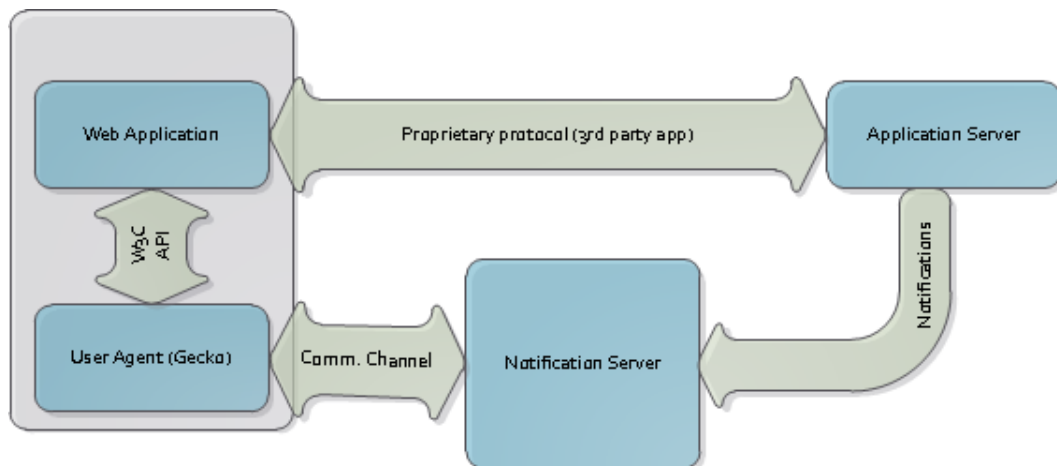
En definitiva, Thialfi puede considerarse como un posible candidato a complementar y no a sustituir a GCM, por las desventajas ya explicadas. La necesidad de crear una segunda conexión y los posibles problemas que puede acarrear a los servidores de aplicación trabajar con millones de peticiones, hace que no sea escalable para pequeñas empresas y malo para dispositivos móviles. Sin embargo, en la práctica se ha demostrado que es competente, aunque con grandes granjas de servidores como los centros de datos de Google.

Capítulo 3. API

El servidor de notificaciones tiene diferentes API (Application Program Interface) que expone hacia los elementos que hay en la arquitectura y que hacen que los mensajes se entreguen a los destinatarios correctos. En un primer momento, y para mantener la compatibilidad con el estándar [<http://dvcs.w3.org/hg/push/raw-file/default/index.html>] que propone el W3C se mantienen los métodos, pero sin embargo, se añaden algunos parámetros extras o se eliminan algunos métodos para añadir nuevas características, seguridad y evitar redundancia.

Pero para entender este capítulo, es necesario presentar a las diferentes instancias que forman parte del servidor de notificaciones:

- **Aplicación web (WebApp: WA).** La aplicación del usuario que es ejecutada en el dispositivo.
- **Agente de usuario (User-Agent: UA).** Como el proyecto se inició bajo el paraguas de Firefox OS, cuyo motor de renderizado de páginas web, equivalente a aplicaciones web es conocido como el Agente de usuario, que en nuestro caso es Gecko, el motor de Mozilla.
- **Servidor de notificaciones (Notification Server: NS).** Es la infraestructura centralizada del servidor de notificaciones. Cualquiera puede desplegar una nueva instancia debido a que está liberado como código abierto [https://github.com/telefonicaid/notification_server]. Como el proyecto se inició bajo el paraguas de Firefox OS, cuyo motor de renderizado de páginas web, equivalente a aplicaciones web es conocido como el Agente de usuario, que en nuestro caso es Gecko, el motor de Mozilla.
- **Servidor de aplicaciones (Application Server: AS).** Es la parte servidora de la WA. Normalmente se podría definir como la presencia en Internet de la WA, que es la que recibe las URLs a las que tiene que hacer push y la que los realiza.

Figura 3.1. Instancias y comunicaciones

Así pues, se dividen la API en dos grandes grupos:

- **Externas.** Son aquellas que se están expuestas en los extremos del servidor, refiriéndose a la comunicación que tienen que realizar las aplicaciones (WA) con el agente de usuario (UA) y el servidor de terceros (AS) hacia el servidor de notificaciones (NS). Son las usadas por desarrolladores externos para que su aplicación sepa recibir notificaciones y su servidor pueda enviarlas.
- **Internas.** Son las que los desarrolladores de aplicaciones no tienen que usar, ya que son transparentes para ellos y están creadas para que el protocolo funcione correctamente, llevando mensajes de señalización y de registro. Están involucrados principalmente el agente de usuario (UA) y el servidor de notificaciones central (NS).

Externas

Con estas APIs, los desarrolladores son capaces de hacer que los mensajes que quieran enviar desde un servidor externo lleguen correctamente hacia sus aplicaciones. La interfaz y el código es muy sencillo, para hacer que el desarrollo sea muy rápido.



Nota

Este API aún no está estandarizado por el W3C, por lo que lo aquí explicado es una proposición funcional.

Hay dos APIs principales externas expuestas para los desarrolladores de aplicaciones, que son:

Entre la WA y UA

Este API, llamada `navigator.mozPush`, es la principal vía de comunicación entre las aplicaciones y el agente de usuario, o sistema operativo.



Nota

Este API está prefijada en la primera versión funcional con `moz` puesto que se ha implementado por primera vez en un dispositivo por parte del navegador Mozilla Firefox y es el que tienen establecido para estas funcionalidades experimentales.

Está basada en el estándar propuesto¹ por el W3C, pero sin embargo tiene una serie de adiciones para añadir seguridad y eliminar la redundancia de alguno de sus métodos.

Con esta API, la aplicación es capaz de realizar dos funciones principales:

- Requerir permiso remoto para poder recibir una URI y comprobar si ya se tiene permiso, mediante la función `requestRemotePermission`.
- Revocar el recurso URI, para eliminar el registro de determinada aplicación, y no recibir más mensajes push, usando `revokeRemotePermission`

Así pues, la interfaz que define a `PushManager` (que es el contenedor que se expone a las aplicaciones), quedará de la forma siguiente:

Figura 3.2. Interfaz `PushManager`

```
interface PushManager {
    PushService requestRemotePermission (DOMString waToken,
                                         DOMString publicKey,
                                         optional DOMString algorithm);
    PushService revokeRemotePermission ();
};
```

Método `requestRemotePermission`

Este método es el que tienen que llamar las aplicaciones para pedir permiso remoto para poder recibir las notificaciones. Así pues, es una petición que tiene dos resultados posibles, `onsuccess` y `onerror`, dependiendo de si el resultado es satisfactorio o no. Como se he especificado más arriba, la signatura de la función es:

Figura 3.3. Función `requestRemotePermission`

```
requestRemotePermission (DOMString waToken,
                        DOMString publicKey,
                        optional DOMString algorithm);
```

Donde se ven los diferentes parámetros necesarios para poder operar correctamente, como son:

¹<http://dvcs.w3.org/hg/push/raw-file/default/index.html>

- **waToken.** Es el token único que tiene que generar la aplicación (ya sea en el servidor, o bien totalmente aleatorio) para identificarse de forma única junto a la clave pública en el servidor de notificaciones, que serán usados para crear una URL unívoca que servirá para mandar las notificaciones solamente a aquellas aplicaciones que hayan registrado ese par.

**Nota**

El desarrollador tiene la opción de hacer su `waToken` público y compartido, por lo que será usado para hacer broadcast de las notificaciones, o bien mantenerlo como un secreto para que sólo la instancia que contenga este token único sea la que reciba las notificaciones.

- **publicKey.** Es una llave pública que la aplicación tiene que tener guardada en el momento de su instalación (o bien transferida de forma segura por el servidor de terceros (AS: Application Server)) para poder ser usada en la comprobación de la firma de los mensajes por parte del componente `NS_AS` (Notification Server - Application Server) del servidor de notificaciones, para prevenir que en el sistema entren notificaciones y mensajes inválidos o maliciosos. Esta llave debe ser transmitida en Base64.
- **algorithm.** [Opcional]. Este parámetro indica el método de firmado usado por el parámetro `publicKey` para la comprobación de las notificaciones. El algoritmo predeterminado es el `RSA-SHA256`.

Además, los diferentes métodos serán llamados una vez se ejecute correctamente o con algún error la función, bien por un problema en la red, o bien por un problema en el servidor:

- **onsuccess.** Este método es invocado cuando la petición `requestRemotePermission` se ejecuta de forma satisfactoria. Lo que realmente hace es recoger un evento `success` que lanza el objeto `requestRemotePermission`. Así entonces, este evento tiene, como en la mayoría de los creados por objetos DOM en los navegadores, un atributo `detail` que contiene el mensaje satisfactorio:
 - **event.detail.url.** URL recibida desde el servidor de notificaciones y que será, probablemente, única para la instancia actual de la aplicación. Esta URL debe ser notificada por el medio más conveniente por la aplicación hacia su servidor de terceros, ya que será la usada para poder recibir notificaciones push.
- **onerror.** Método invocado cuando sucede algo no esperado en la respuesta del objeto creado por `requestRemotePermission`. Puede ser debido por múltiples factores, como fallo en la red (no disponible, o que se haya caído en medio de la transmisión), fallo en el servidor de

notificaciones (mensaje de error por no disponibilidad), o incluso por que el registro falle por tener un token incorrecto o duplicado.

El evento que se captura también tiene el mensaje de error, que si bien no es informativo para las aplicaciones (debería ser completamente transparente), se muestra en:

- **event.error.** Mensaje de error. Bien generado por el propio sistema operativo o UA (User Agent) o bien el que es enviado por el servidor de notificaciones.

Ejemplo 3.1. Pedir permiso remoto y recibir URL

```
var watoken = 'mySecretWaToken';
var pbk = '...';

var push = navigator.mozPush;
var req = push.requestRemotePermission(watoken,
                                       encodeBase64(pbk));

req.onsuccess = function(event) {
    console.log('Yay! URL is: ' + event.detail.url);
};

req.onerror = function(event) {
    console.log('Oops!! Error is: ' + event.error);
};
```

Además, hay otro método, función o callback, llamado `onmessage`, que es ejecutado cada vez que la aplicación recibe un mensaje de notificación. Sólo es ejecutado cuando ha habido un `onsuccess` previo, y no se ha producido ningún `onerror`.

- **onmessage.** Función que se llama una vez llega una notificación desde el sistema operativo. Hay que notar, que para que funcione esta opción, el objeto creado por `requestRemotePermission` tiene que seguir vivo y no haber sido borrado por el recolector de basura de JavaScript, por lo que se puede interpretar con que la aplicación está abierta, o bien se crea un objeto con una referencia fuerte, que significa que si la aplicación se cierra, el objeto sigue abierto y puede ser notificado.
- **event.detail.message.** Mensaje completo. Es el `payload` que el servidor de las aplicaciones quería mandar desde un primer momento.
- **event.detail.url.** [Opcional] URL a la que la notificación iba dirigida. Añadido por la posibilidad de que en un futuro las aplicaciones puedan pedir varias veces el permiso remoto, por lo cual recibirán diferentes URL a la que el servidor podría hacer push.



Nota

En la primera versión implementada en FirefoxOS, se utiliza el futuro estándar `setMessageHandler` del objeto `navigator` para poder ser despertado a pesar de que la aplicación esté cerrada.

```
navigator.mozSetMessageHandler("push-notification", function(msg) {  
    alert("New Message with body: " + JSON.stringify(msg));  
});
```

Método `revokeRemotePermission`

Figura 3.4. Función `revokeRemotePermission`

```
revokeRemotePermission ();
```

Este método realiza el paso contrario que el anterior, puesto que pide al servidor de notificaciones que se dejen de enviar notificaciones, esto es, que se quiere eliminar la URL o recurso dado en el paso anterior. Es un de-registro de URL. Puede servir en múltiples casos, como por ejemplo:

- Si la aplicación decide eliminar las notificaciones push por otro sistema y quiere dejar todo el sistema limpio, haciendo olvidar y eliminar todos los datos anteriores.
- El usuario decide editar las preferencias de la aplicación desde la configuración global del sistema y desactiva la opción de notificaciones push para determinadas aplicaciones. El sistema puede decidir no mandar las notificaciones a la aplicación o bien de-registrarla del servidor.
- El servidor de terceros ha sido comprometido y las aplicaciones tienen que re-registrarse para usar la nueva clave pública (P_{bK}), por lo que en primer lugar tienen que de-registrarse para evitar que futuras notificaciones de orígenes malignos puedan volver a llegar a la aplicación.

Ejemplo 3.2. Pedir revocación de permiso (revocar URL)

```
var push = navigator.mozPush;
var req = push.revokeRemotePermission();

req.onsuccess = function(event) {
    console.log('Unregister successful');
};

req.onerror = function(event) {
    console.log('Oops!! Error is: ' + event.error);
};
```

Entre el WA y AS

La comunicación entre la aplicación web (WA) y el servidor de terceros (AS: Application Server) está fuera del alcance del servidor de notificaciones, ya que depende de cada aplicación el cómo se implementa.

Sin embargo, sí hay una serie de medidas que se pueden tomar para que esta transmisión de un dato importante sea seguro y no pueda ser interceptado por terceros que, aunque después sea complicado su uso porque depende de poder firmar con la clave privada las notificaciones, puede llevar a conocer datos privados de los clientes o instancias si se conoce el método de generación del `appToken` (que está contenido en la URL) y la clave pública.

- Que el envío se realice de forma segura, mediante HTTPS, para que sólo los extremos puedan conocer el contenido de la comunicación.
- Usar cifrado para evitar que un posible ataque de man-in-the-middle pueda recoger los datos y saber de qué se trata. Algo sencillo computacionalmente ya que el cifrado se lleva a cabo en el cliente, que suelen tener menor potencia de cálculo que los servidores a los que van dirigidos.
- No pensar sólo en el envío HTTP, sino que otros sistemas podrían ser correctos, como XMPP, correo electrónico o incluso un mensaje SMS si fuera necesario.
- Por supuesto, asociar algún número que no permita saber a quién corresponde la URL enviada con qué usuario para un atacante externo, pero que el servicio pueda identificarlo de forma unívoca.

Entre el AS y NS

Este es el API que tienen que utilizar los servidores web para enviar las notificaciones de una forma estándar hacia el servidor de notificaciones para que éste las mande correctamente hacia los clientes.

El API se basa en protocolos estándar de la web y bien conocidos, como son los HTTP Post, mensajes JSON (JavaScript Object Notation) y las cabeceras de respuesta HTTP estándar, por lo que es muy sencilla la implementación y el protocolo, a parte de ser muy ligero y de haber cientos de programas software o librerías que implementan los métodos necesarios para poder enviar estas notificaciones.

Así pues, el servidor de terceros (AS: Application Server), necesita varios elementos antes de poder enviar su notificación, como son:

- Una clave privada, para poder firmar las notificaciones. Tiene que ser la correspondiente a la clave privada que tienen las aplicaciones guardada y con la que se han registrado.
- Una URL (o recurso unificado) a la que enviar las notificaciones. Habrá sido enviada por la aplicación en un paso anterior por algún método elegido por el desarrollador.
- Un mensaje que enviar (ya sea un número de versión o un mensaje con poco contenido o hasta 4KiB).

Entonces, una vez que el servidor de terceros tenga estos datos, es capaz de empezar a mandar notificaciones para que sean entregadas a las aplicaciones finales, siempre dependiendo de la URL y de la firma con la que se mandan dichas notificaciones.

**Nota**

Este API se corresponde a una simple interfaz REST que será mejorada en futuras versiones para incluir más características.

Así pues, se acepta una petición HTTP Post usada para mandar la notificación con el mensaje real dentro. El siguiente esquema debe ser enviado la `publicURL` que defina a la instancia de la aplicación para un usuario determinado (explicado en los puntos anteriores), del estilo de: `https://push.telefonica.es/notify/SOME_RANDOM_TOKEN`

Figura 3.5. Notificación enviada

```
{
  messageType: "notification",
  id: "<ID used by the Application Server>",
  message: "<Message payload>",
  signature: "<Signed message>",
  ttl: "<time to live>",
  timestamp: "<Since EPOCH Time>",
  priority: "<prio>"
}
```

Así, como se puede ver la notificación tiene las siguientes características:

- Es una notificación con un estilo JSON, simple y sencillo, que es autocontenido y fácil de parsear y de crear.
- El tamaño de la notificación es mínimo, menor que si el sistema de intercambio de mensajes fuera un estricto XML.

Y define los siguientes atributos (obligatorios si no se indica lo contrario):

- **messageType.** Tiene que ser igual al texto `notification`. En un futuro podría haber más `messageType` para pedir información sobre la entrega de los mensajes o más funcionalidades.
- **id.** Identificador interno para el servidor de terceros, tiene que ser una cadena menor de 32 caracteres. Actualmente no tiene uso, pero en un futuro servirá para poder sobrescribir notificaciones (mandar otra notificación con el mismo `id` que el anterior y, si aún no se ha entregado, modificarla) y para poder preguntar cuál es el estado de la notificación: entregada, sin entregar, expirada...
- **message.** Es el texto que tiene que utilizar la aplicación y que el servidor de terceros envía. Tiene que ser un texto plano o bien estar en un formato JSON, pero convertido a cadena, por ejemplo, con la función `JSON.stringify` de JavaScript.
- **signature.** La firma con la clave privada de lo que se manda dentro de `message`. Se firma con `RSA-SHA256` por defecto, o con el algoritmo que se haya definido en el atributo `algorithm` pasado a la función `requestRemotePermission` cuando la aplicación ha pedido permiso a su sistema operativo.
- **ttl.** [Opcional] Tiempo de vida del mensaje. Está especificado en segundos, y puede ser desde 0, lo cual significaría que no se guarda en la base de datos, y si los destinatarios están desconectados, o no se puede entregar, la notificación se pierde; hasta 4 semanas (o 31 días, 2678400 segundos). El tiempo por defecto es de 4 semanas. Una vez superado ese tiempo, si aún sigue en el sistema (esto es, no se ha entregado), se elimina y no se entrega.
- **timestamp.** [Opcional] Tiempo desde el inicio UNIX, en segundos, para indicar cuándo salió la notificación desde el servidor de terceros. A esto se le suma el `ttl` para el tiempo máximo de vida en el sistema. Por defecto, se pone el momento en el que se recibe, si no está indicado.
- **priority.** [Opcional] Prioridad del mensaje. Varía entre 1 y 4, donde 1 es más prioritario, y 4 es menos prioritario. Por defecto el valor que se indica es 4. Usado para mostrar al servidor si tiene que entregar la notificación en el momento o puede retrasarla hasta que haya más o el terminal esté despierto.

Pero además, se definen una serie de respuestas por parte del servidor dependiendo de si las notificaciones han sido aceptadas o no:

Mensajes sin códigos de error

Sólo hay un mensaje de OK, que indica que la notificación se ha procesado correctamente y va a poderse entregar a los dispositivos que están vinculados a la URL a la que se ha hecho HTTP Post.

```
Estado HTTP: 200
{
  status: "ACCEPTED"
}
```

Mensajes de error del servidor

Pero los mensajes más interesantes son los de error, que son todos aquellos que sean diferentes al de la cabecera HTTP con respuesta 200. Dependiendo de la cabecera HTTP, el error varía, dando la razón en un formato escrito en el atributo `reason` de la respuesta, que es un JSON.



Nota

Estos mensajes de error están definidos en el fichero `/src/common/constants.js` a partir de la línea 28²

- **JSON inválido.** Se produce cuando la notificación no está correctamente formada, y no puede ser parseada como un JSON válido. Puede que sea debido a que no hay un atributo correctamente cerrado o contenga caracteres no parseables dentro del juego de caracteres de JavaScript (UTF-16).

```
Estado HTTP: 450
{
  status: "ERROR",
  reason: "JSON not valid error"
}
```

- **URL incorrecta.** No se ha reconocido el formato de la URL y no se puede identificar a quién va dirigido.

```
Estado HTTP: 404
{
  status: "ERROR",
  reason: "Bad URL"
}
```

²https://github.com/telefonicaid/notification_server/blob/develop/src/common/constants.js#L28

- **appToken inválido.** El appToken enviado no es válido o está vacío.

```
Estado HTTP: 451
{
  status: "ERROR",
  reason: "No valid appToken"
}
```

- **Se está usando un método HTTP no soportado.** No se está enviando una petición HTTP Post.

```
Estado HTTP: 405
{
  status: "ERROR",
  reason: "No valid HTTP method"
}
```

- **No es una notificación.** El atributo messageType no contiene la cadena notification, por lo que no es válido.

```
Estado HTTP: 452
{
  status: "ERROR",
  reason: "Not messageType=notification"
}
```

- **La notificación no está firmada.** La notificación no tiene una firma para poder verificar el origen.

```
Estado HTTP: 453
{
  status: "ERROR",
  reason: "Not signed"
}
```

- **Identificador no válido.** El identificador, que tiene que ser una cadena no nula, es inválido.

```
Estado HTTP: 454
{
  status: "ERROR",
  reason: "Bad id"
}
```

- **La notificación es demasiado grande.** Supera el límite del servidor de notificaciones, que está en 4KiB ó 4096 bytes.

```
Estado HTTP: 413
{
  status: "ERROR",
  reason: "Body too big"
}
```

- **Firma incorrecta.** La notificación está firmada, pero no es correctamente válida, esto es, la comprobación de la firma con la clave pública con la que la aplicación se registró ha devuelto que el mensaje ha sido modificado o es inválido.

```
Estado HTTP: 455
{
  status: "ERROR",
  reason: "Bad signature, dropping notification"
}
```

A partir de estos errores, el servidor de terceros es capaz de determinar que problema tiene el servidor de notificaciones con las notificaciones que está intentando enviar.



Nota

Los mensajes de error pueden cambiar de un servidor a otro, sobre todo si se añaden extensiones o más características, pero estos son los básicos que debería responder en caso de algún problema.

Así pues, las secciones anteriores son las APIs externas con las que tendrán que trabajar los desarrolladores de aplicaciones y gestores de servidores que intentan enviar notificaciones push.

Interna: Entre el UA y el NS

Estas APIs y comunicaciones son las que se producen internamente, de forma transparente a los desarrolladores de las aplicaciones y del servidor que envía la notificación. Así pues, tiene que ser tratada por los desarrolladores que quieran implementar el protocolo en su plataforma, o sistema operativo.

Está dividida en tres partes principales, una interfaz HTTP, para recoger el token, otra que usa el protocolo WebSockets, para el intercambio bidireccional de datos y una tercera, que es la interfaz UDP que tiene que tener el teléfono, así pues:

Interfaz HTTP

A través de este API, se puede descargar el token inicial y único que será usado para el registro posterior y demás operaciones.



Nota

Este API puede cambiar en un futuro y entre diferentes servidores de notificaciones, puesto que el token podría venir predeterminado en el mismo dispositivo o ser recogido mediante otros mecanismos, como correo electrónico, respuesta a un registro web, etcétera.

Así pues, la petición actual es:

```
Host: https://push.telefonica.com/  
GET /token
```

Y la respuesta sería:

```
Status: 200  
Body: a4b26ecbd961c673f3526a8cc747758...
```

Donde todo el cuerpo de la respuesta (`res.body`) sería el token de dispositivo que es el que se tendría que usar para realizar los registros de dispositivo y de aplicaciones.

Interfaz WebSocket

Este es el API bidireccional para el intercambio de datos, registros y notificaciones en algunos casos. Es el más complejo de todos ya que están involucrados muchos mensajes diferentes, pero todos son correctamente formados como JSON.

En primer lugar, es necesario abrir una conexión WebSockets con el servidor de notificaciones (en concreto con la instancia `NS_UA_WS`) para poder empezar a intercambiar datos.

Conexión WebSocket

Así pues, hay que utilizar el recurso necesario para iniciar la conexión, del estilo de `wss://push.telefonica.com`, usando el protocolo `push-notification`.



Nota

A partir de este momento, los siguientes métodos y comunicaciones se realizan por el canal WebSocket que se ha abierto en la sección anterior.

Registro de UA

Es el primer paso que hay que realizar para empezar a recibir notificaciones. Este paso significa el registro del dispositivo en el servidor de notificaciones, y la manera de recibirlas. Se realiza enviando un mensaje al servidor por la conexión WebSocket de la manera:

```
{
  messageType: "registerUA",
  data: {
    uatoken: "<a valid UAToken>",
    interface: {
      ip: "<current device IP address>",
      port: "<TCP or UDP port to wake up>"
    },
    mobilenetwork: {
      mcc: "<Mobile Country Code>",
      mnc: "<Mobile Network Code>"
    }
  }
}
```

Donde los diferentes atributos tienen que ser:

- **messageType.** Tiene que ser siempre `registerUA`.
- **data.** Otros datos, compuestos por `uatoken`, `interface` y `mobilenetwork`.
 - **uatoken.** Es el token de dispositivo conseguido en las secciones anteriores.
 - **interface.** Contiene la IP y puerto en el que el dispositivo puede ser despertado. Por ahora, sólo se permite la escucha en puertos TCP o UDP.
 - **ip.** Dirección del dispositivo, ya sea con la IP privada o pública. En el caso de que sea pública, no hace falta que se rellene el dato de `mobilenetwork` y en el caso de privada, sí es necesario, puesto que se usará un proxy intermedio colocado dentro de la red móvil para poder encontrarle directamente.
 - **port.** Puerto de escucha, puede ser TCP o UDP. No necesitará saber qué le llega, sino simplemente que algo ha llegado ahí.
 - **mobilenetwork.** Indica los dos valores necesarios para identificar una red móvil correctamente, `mcc` y `mnc`. Las listas completas pueden consultarse en Internet³.

³http://en.wikipedia.org/wiki/Mobile_country_code

- **mcc.** Mobile Country Code.
- **mnc.** Mobile Network Code.

Así pues, la respuesta del registro puede ser válida o no.



Nota

Hay que notar que las respuestas son asíncronas, y el envío de un mensaje de tipo `registerUA` no significa que su mensaje de respuesta sea el siguiente que llegue, por lo que hay que tratar cada mensaje por separado sin esperar una respuesta instantánea.

En el caso de una respuesta exitosa, se recibirá algo similar a:

```
{
  "messageType": "registerUA",
  "status": "REGISTERED",
  "pushMode": "WS",
  "WATokens": [],
  "messages": []
}
```

Donde se pueden encontrar varios atributos diferentes, los cuales son:

- **messageType.** Que tiene que ser igual a `registerUA`, ya que es la respuesta al tipo de mensaje enviado en un inicio.
- **status.** Si el valor es `REGISTERED`, es que el registro ha sido correcto.
- **pushMode.** Indica cuál es el modo en que las notificaciones van a llegar al dispositivo. Hay varios modos, los cuales pueden ser
 - **WS.** WebSocket. La conexión debe mantenerse abierta para poder recibir las notificaciones.
 - **UDP.** La conexión puede ser cerrada, ya que el dispositivo va a ser informado de que tiene nuevas notificaciones usando el canal UDP y el puerto abierto en el dispositivo. Útil en redes móviles.
 - **wapPush.** Utilizará WAP Push para informar al dispositivo de que hay nuevas notificaciones. No implementado en la primera versión



Nota

No están contempladas todas las posibilidades, ya que se puede añadir más en el futuro, dependiendo de los dispositivos que estén funcionando.

- **WATokens.** Array que indica todas las aplicaciones que están registradas en el dispositivo. Se mandan las URLs completas, por lo que se puede comprobar si falta alguna. En el caso de que el dispositivo tenga más de las recibidas, puede re-registrar la aplicación determinada o bien eliminarla, según el caso. Si no ha recibido una aplicación que debería estar registrada, simplemente deberá registrarla.
- **messages.** Array con las diferentes notificaciones que están pendientes de recibir por el dispositivo. Por cada una de ellas, y si el procesado es correcto, se debe enviar un mensaje de confirmación, o `ACK` para indicar que se ha recibido correctamente y no enviar duplicados en un futuro.

En el caso de una respuesta incorrecta, la respuesta será⁴:

```
{
  "messageType": "registerUA",
  "status": "ERROR",
  "reason": "<razón>"
}
```

Donde la razón `<razón>` puede ser:

- **Data received is not a valid JSON package.** El mensaje enviado no es válido.
- **UAtoken not valid for this server. Get a new one.** El token no es válido para el servidor actual, lo que quiere decir que se rechazará hasta que no se consiga uno nuevo.
- **Failed registering UAtoken.** Hubo un fallo interno en el servidor al registrar el dispositivo.



Nota

Esta función de registro tiene que realizarse cada vez que cambie la conexión móvil o WiFi del dispositivo, para actualizar la información en el servidor.



Nota

A partir de este momento, los siguientes métodos y comunicaciones necesitan tener una conexión autorizada, esto es, que haya registrado anteriormente un dispositivo mediante el método `registerUA`

⁴https://github.com/telefonicaid/notification_server/blob/develop/src/common/constants.js#L39

Registro de WA

Esta función es la usada para el registro de aplicaciones a partir de una clave pública y un token de aplicación. Es la capa inferior de la función `requestRemotePermission` del objeto `navigator.push` expuesto para los desarrolladores de aplicaciones.

El formato es como sigue:

```
{
  "data": {
    "watoken": "hola",
    "pbkbase64": "hola"
  },
  "messageType": "registerWA"
}
```

Donde los diferentes atributos se definen como:

- **data.** Contiene los datos necesarios para realizar el registro, los cuales son:
 - **watoken.** Es el token de aplicación que hay que registrar.
 - **pbkbase64.** Es la clave pública, codificada en formato Base64 para evitar problemas con los saltos de línea u otro tipo de caracteres o símbolos extraños, que corresponde a la aplicación que ha solicitado el registro.
- **messageType.** Tiene que ser igual a la cadena `registerWA`.

Con la posible respuesta, por el mismo canal y de forma asíncrona que puede ser:

```
{
  "watoken": "<watoken>",
  "messageType": "registerWA",
  "status": "REGISTERED",
  "url": "https://push.telefonica.com/notify/5819b005..."
}
```

Que se corresponden con los siguientes atributos:

- **watoken.** Es el token que ha generado la respuesta actual. De nuevo, comentar que las respuestas son asíncronas, por lo que tiene que ir reflejado a qué petición inicial iba dirigida.
- **messageType.** Tiene que ser igual al mensaje original, por lo tanto `registerWA`. Permite tratar al mensaje correctamente en el código.

- **status.** Tiene que ser igual a la cadena `REGISTERED`. En caso contrario, o que venga el texto `ERROR`, el registro no ha sido correcto.
- **url.** URL o recurso con la dirección única de notificación. Tiene que enviarse esta URL al objeto que ha originado la petición desde la aplicación.

También hay unos códigos de error similares a los del método `registerUA` respondiendo con un mensaje del estilo:

```
{
  "watoken": "<watoken>",
  "messageType": "registerWA",
  "status": "ERROR",
  "reason": "<razón>"
}
```

Con el atributo `reason` dando una explicación del error en forma humana, con un claro mensaje descriptivo y un atributo `watoken` indicando qué aplicación ha fallado.

Una vez recibido un mensaje correcto de registro de WA, se puede estar seguro de que las notificaciones que envíen los servidores de terceros eventualmente llegarán a la aplicación a la que se envía la notificación.

Recepción de notificación

Las notificaciones pueden venir por dos medios diferentes

- **En respuesta a `registerUA`.** Cuando un dispositivo se registra o actualiza sus datos, tiene que enviar un mensaje del tipo `registerUA`, el cual, si todo funciona correctamente, recibirá de forma asíncrona la respuesta, conteniendo un array de mensajes en el atributo `messages`.

```
{
  "messageType": "registerUA",
  "status": "REGISTERED",
  "pushMode": "WS",
  "WATokens": [],
  "messages": []
}
```

- **Como mensaje independiente.** Sin embargo, si se está utilizando una conexión de WebSockets, las notificaciones llegan como mensajes independientes, de la manera siguiente:

```
{
  "messageType": "notification",
  "id": 1234,
  "message": "Hola",
  "ttl": 2592000,
  "timestamp": "1359984038",
  "priority": 1,
  "messageId": "aacf9ca0-6ecd-11e2-815f-551df8e4ba6a",
  "url": "https://push.telefonica.com/notify/31387b27716ac..."
}
```

Donde lo más importante son los dos últimos atributos:

- **messageId.** Es el identificador único que le ha dado el servidor de notificaciones a la notificación. Es el usado para mandar el acuse de recepción.
- **url.** URL inicial a la que se envió la notificación, sirve para identificar en el User Agent (UA) la aplicación en concreto a la que va dirigida.

El UA que recibe este mensaje debe utilizar el atributo `url` para saber a qué aplicación en concreto va, y una vez encontrada. Debe eliminarse de la notificación entregada a la aplicación.

El UA también debe usar el atributo `messageId` para enviar la confirmación de recepción, una vez esté recibido correctamente por la aplicación. Debe eliminarse de la notificación entregada a la aplicación.



Nota

Para no volver a recibir los mensajes duplicados, hay que enviar un acuse de recibo, explicado en la siguiente sección.

Confirmación de recepción, o ACK

Una vez recibidas notificaciones, ya sea a través de los mensajes de registro de User Agent (a través del mensaje `registerUA`) o mediante la conexión WebSocket, una notificación directa, hay que mandar un mensaje de recepción al servidor, por varios motivos:

- Al aceptar el mensaje, éste se borra de la base de datos, en concreto de la colección de mensajes dirigidos para determinado usuario, por lo que si no se informa de su recepción, el mensaje seguirá enviándose hasta que pase el tiempo de vida (especificado en el atributo `ttl`) o se de-registre la aplicación.
- Los servicios que envíen la notificación, serán capaces de controlar si la notificación se ha entregado correctamente o con algún error a los diferentes dispositivos que la deberían recibir.

Así pues, el mensaje que hay que enviar a través de la conexión WebSocket para informar de la recepción de los mensajes es el siguiente:

```
{
  messageType: 'ack',
  messageId: <msg.messageId>
}
```

Donde los diferentes atributos tienen que ser:

- **messageType.** Tiene que ser igual a la cadena `ack`.
- **messageId.** Tiene que ser igual al `messageId` que está contenido dentro del mensaje. No es el `id` que ha enviado el servidor de terceros, sino uno generado internamente para el guardado y búsqueda del mensaje en el servidor de notificaciones.

No hay ninguna respuesta a este mensaje y además tiene que enviarse uno por cada mensaje que haya que aceptar la recepción. Así pues, si en la respuesta a `registerUA` llega un array con 10 mensajes diferentes, hay que mandar 10 mensajes de recepción del estilo superior para que se borren de la base de datos.

De-registro de WA

De la misma manera que se puede registrar una aplicación, también puede de-registrarse para que el dispositivo no siga recibiendo notificaciones continuamente en los siguientes casos:

- El usuario ha desinstalado la aplicación, por lo que el gestor de aplicaciones, al eliminarla, también tiene que eliminar los permisos asociados y realizar las operaciones oportunas. En el caso de push, lo que se tiene que hacer es de-registrar la aplicación en local y en remoto, para que la URL no sea válida y pueda dejar de recibir notificaciones push.
- El usuario revoca el permiso en la aplicación de configuración, donde se le presenta una lista de aplicaciones que están utilizando el servicio de push.
- La aplicación pide la revocación programáticamente, porque no quiere seguir recibiendo nuevos mensajes push, o bien por algún cambio de usuario o de configuración que requiera que sea otra URL la de recepción push.

Entonces, el mensaje que se tiene que enviar es el siguiente:

```
{
  "data": {
    "watoken": "<watoken>",
    "pbkbase64": "<pbkbase64>"
  },
  "messageType": "unregisterWA"
}
```

Donde los atributos que hay que enviar son los siguientes:

- **data.** Que contiene los datos siguientes:
 - **watoken.** Que es el token de aplicación dado desde la aplicación que solicita el de-registro. Tiene que ser una cadena alfanumérica.
 - **pbkbase64.** La clave pública para el de-registro, tiene que estar codificada en base64 para evitar problemas en la transmisión con saltos de línea y otros símbolos no serializables.
- **messageType.** Tiene que ser igual a la cadena `unregisterWA`.

Y la respuesta, de nuevo asíncrona, es como sigue:

```
{
  "url": "https://push.telefonica.com/notify/31387b27716ac...",
  "messageType": "unregisterWA",
  "status": "UNREGISTERED"
}
```

Donde se indica, como en los mensajes anteriores, diferentes atributos:

- **url.** Que se corresponde a la URL que ha sido de-registrada.
- **messageType.** Tiene que ser igual a la cadena `unregisterWA` para especificar que es la respuesta a una petición similar anterior.
- **status.** Si todo ha sido correcto, el atributo `status` contendrá una cadena igual a `UNREGISTERED`

De-registro de UA

De la misma manera que hay un método de registro de UA y un método de de-registro de WA, hay un método para de-registrar un dispositivo.



Nota

Peligro. Este método elimina todos los datos del dispositivo de la base de datos, invalidando cualquier URL y desechando cualquier notificación pendiente.

El mensaje que hay que enviar, es simplemente:

```
{
  "messageType": "unregisterUA"
}
```

El cual es el único método síncrono, ya que si el de-registro es correcto, la conexión WebSocket por la que se ha enviado, se cierra de forma consensuada (esto es, el cierre viene acompañado por el código 1000 de las constantes de cierre de WebSocket)⁵.

El mensaje de cierre, aunque sea nulo, podría ser en la consola JavaScript algo similar a:

```
Closed connection to push.telefonica.com with code 1000
and reason [undefined]
```

A partir de este momento, la conexión está cerrada, y el dispositivo no existe para el servidor de notificaciones. Para poder volver a registrarse, se tienen que repetir los pasos superiores.

Pedir WA registradas

Este es un método para poder mantener una sincronía y el mismo estado entre el servidor y el cliente. La lógica de qué hacer con los datos recibidos la tiene que tener el cliente, por ejemplo:

- Si vienen más URLs de las que están en local, quizás haya que de-registrar alguna aplicación. Pero también puede ser que se haya perdido el registro local.
- Si vienen menos URLs de las que debería haber, habría que registrar de nuevo una aplicación, puesto que el servidor no sabe nada de ella y no podría recibir notificaciones push.

La petición se realiza por el canal abierto, y el mensaje a enviar es el siguiente:

```
{
  "messageType": "getRegisteredWA"
}
```

Con una posible respuesta como esta:

⁵<http://tools.ietf.org/html/rfc6455#section-7.4.1>

```
{
  "WATokens": [
    "https://push.telefonica.com/notify/31387b27716ac...",
    "https://push.telefonica.com/notify/739a6ecd89f3e...",
    "https://push.telefonica.com/notify/113ef443accd3..."
  ],
  "messageType": "getRegisteredWA"
}
```

Donde los atributos son los siguientes:

- **WATokens.** Es un array con todas las URLs de aplicaciones registradas.
- **messageType.** Tiene que ser igual a la cadena `getRegisteredWA` puesto que es la respuesta a la petición anterior.

Interfaz UDP

Es una interfaz de un único sentido, en el cual el servidor de notificaciones informa a un dispositivo, o UA que tiene una notificación, y que debe conectarse al servidor para recogerla.

Es la parte más diferenciadora de este servidor de notificaciones con respecto al resto de opciones que hay ahora mismo en otros sistemas operativos, porque es la responsable de que el uso de batería sea menor y el número de puertos abiertos en el GGSN debido a aplicaciones intentando mantener sockets abiertos, disminuya drásticamente.

Así pues, este interfaz lo único que requiere es que el dispositivo tenga un puerto abierto de escucha, por el que recibirá pequeños paquetes de datos, o PINGs para indicar que tiene, e indicar en el registro del dispositivo cuál es su IP y su puerto, así como el MCC y MNC de su conexión móvil.

A partir del momento del registro, si todo es correcto, y el servidor tiene un proxy registrado para poder mandar dichos pings a los dispositivos que están en una red privada, la conexión WebSockets se cerrará y se utilizará dicho mecanismo de PING.

Los PINGs son de tipo UDP o TCP, pero nunca van a llevar un contenido que tenga que ser leído por el teléfono, sino que simplemente, al recibir este paquete, tiene que conectarse con el servidor y realizar un registro de nuevo para recibir las notificaciones en la respuesta.

Flujo completo

Para poner todo en orden, estos son los pasos que tiene que realizar el User Agent o sistema operativo para poder utilizar las notificaciones push y que terceros puedan avisarle:

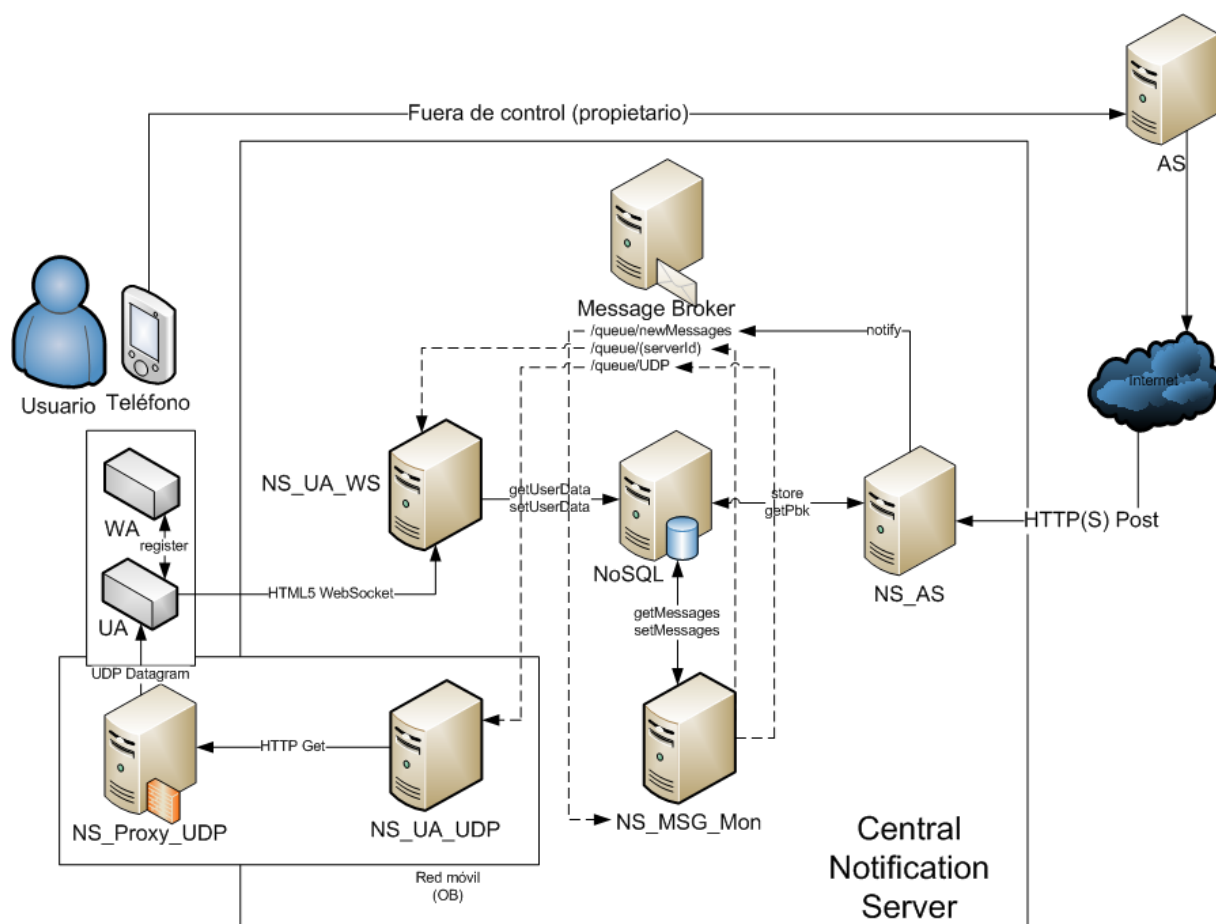
1. Si no tiene un token de dispositivo (`uaToken`), conseguir uno.
2. Abrir una conexión WebSockets con el servidor de notificaciones usando el protocolo `push-notification`.
3. Registrar el dispositivo con el mensaje del tipo `registerUA`
 - a. Si es correcto, procesar los posibles mensajes que vengan en el atributo `messages` y continuar al paso siguiente.
 - b. Si es incorrecto, comprobar la causa, arreglarla, y reintentar.
4. Registrar una aplicación cuando sea pedido por las aplicaciones.
5. Si llega un mensaje de notificación:
 - a. Comprobar si se puede entregar a la aplicación (si aún existe).
 - b. Entregarlo y mandar un mensaje de `ACK` al servidor para que no lo vuelva a mandar.
6. Aceptar y procesar cualquier mensaje de registro de aplicaciones, `registerWA`, mediante el guardado en el propio sistema de cuál es la aplicación que lo pide (guardando alguna referencia a ella, puesto que en un futuro hay que despertarla) y la URL que devuelve el servidor de notificaciones para poder identificar las notificaciones dirigidas a ella.
7. Cada vez que el sistema cambie su conexión, ya sea a red móvil o wifi, tiene que informarse al servidor de dicho cambio, incluyendo `IP`, `port`, `MCC` y `MNC` si corresponde.
8. En caso de de-registro total, eliminar toda la base de datos local y pedir el de-registro del servidor mediante los mensajes `unregisterWA` y `unregisterUA`.
9. Enviar mensajes de `PING` en el caso de que sea necesario o apropiado.

Capítulo 4. Arquitectura interna

El capítulo tiene como finalidad explicar cómo se ha organizado internamente el servidor de notificaciones central y cómo interactúan las diferentes partes o instancias entre ellas para que las notificaciones sean tratadas correctamente, desde el momento en que son recibidas hasta cuando el teléfono va a procesarlas y tiene que dar la confirmación de recepción.

Pero antes de explicar el funcionamiento y los diferentes mensajes que se intercambian, se explicará los diferentes actores que están en funcionamiento repartidos por todo el servidor.

Figura 4.1. Esquema general de la arquitectura



Tipos de servidor

Para poder escalar tanto horizontal como verticalmente sin tener ningún límite, la infraestructura del servidor de notificaciones central ha sido dividida en diferentes instancias o cajas, cada una dedicada a una función en particular, y que a su vez son independientes del resto, comunicándose sólo por mensajería.



Nota

Los nombres de cada caja siguen la nomenclatura: NS-<tipo_de_cliente>.

NS-UA-WS

Este servidor es el encargado de ejercer como frontend para los dispositivos móviles, exponiendo una serie de APIs para que los diferentes agentes de usuario (o UA) puedan utilizarlo. En concreto, posee dos interfaces, una HTTP y otra como WebSockets.

Las dos interfaces son:

- **HTTP.** Sirve para pedir un token válido para el agente de usuario, único, que sirve para identificar a cada dispositivo. Este token puede ser recogido llamando al método `HTTP Get` a la url `https://push.telefonica.com/token`
- **WebSocket.** Diseñado para ejercer como canal de comunicación bidireccional entre el agente de usuario y el servidor de notificaciones. Es la conexión por la cual se realizan los registros y se envían las notificaciones.

Y sus principales cometidos son:

- Guardar en la base de datos de MongoDB todos los nodos registrados y las aplicaciones, además de tener que recogerlos en determinado momento.
- Suscribirse a una cola de mensajes creada específicamente para él, por el cual le llegan las notificaciones de que tiene que enviar un mensaje a un dispositivo.
- Proveer las dos interfaces explicadas anteriormente, que sirven como nexo de comunicación entre los dispositivos y el servidor de notificaciones.

NS-UA-UDP

Este servidor es el responsable de intermediar entre el servidor de notificaciones central y el servidor de wakeup, o `NS_Proxy_UDP` que está desplegado en el interior de cada operadora móvil.

Al igual que el servidor anterior, el `NS-UA-WS`, este también está escuchando una cola determinada, llamada `UDP`, para que cada vez que llegue un mensaje, lo procese, recupere de la base de datos MongoDB la información necesaria y pueda enviar al servidor `NS_Proxy_UDP` un mensaje HTTP para que éste último despierte al teléfono con un pequeño paquete UDP.

NS_Proxy_UDP

Este servidor es un proxy entre el complejo central del servidor de notificaciones y los clientes. Este servicio recibe peticiones a través de un puerto HTTP estándar, y manda datagramas UDP o paquetes TCP a los dispositivos dentro de la red móvil con el propósito de despertarles y que se conecten a la infraestructura central. Esta instancia debe estar instalada dentro de la zona privada de la operadora, y tener visión directa de todos los terminales móviles en dicha red.

NS-AS

Esta instancia, es el front-end para los servidores de aplicaciones. Proporciona una interfaz HTTP Post para poder recibir notificaciones.

La interfaz expone la ruta `/notify`, en la cual una vez recibido, se procesan los mensajes, comprobando si son correctos, si están firmados y son válidos (recogiendo la clave pública desde la base de datos central) y realizando dos pasos diferentes:

1. Introducir el mensaje en la base de datos MongoDB para mantener una persistencia y no perder el mensaje por si ocurre algún error.
2. Enviar una notificación a la cola `newMessages`, que indica que se ha recibido una nueva notificación.

NS-MSG-Monitor

Es el responsable de recibir los mensajes a través de la cola `newMessages`, procesarlos, encontrar a quién están dirigidos, y enviarlo a la cola de los servidores que dan servicio a los diferentes dispositivos de forma correcta. Se podría decir que por esta instancia pasan todos los mensajes que tienen que moverse por dentro del complejo sistema del servidor. Además, monitoriza el estado de los mensajes y los re-encola si fuera necesario, así como eliminarlos en un futuro si sobrepasan su tiempo de vida máximo

Colas de mensajes, o Message Queue

Son varias instancias del sistema de mensajería RabbitMQ que actúa como transportista de mensajes entre todos los diferentes servidores que lo componen.

En primer lugar, su función es la de creación de colas, que son alimentadas por un extremo, y consumidas por el otro. La configuración varía según las colas y cuál sea su cometido, pero por ejemplo, se puede encontrar:

- **/newMessages.** Recibe todos los nuevos mensajes a partir del servidor `NS-AS`. Los lee el monitor `NS-MSG-Monitor`.
- **/UDP.** Una cola que leen todos los servidores de UDP, `NS-UA-UDP` y que es escrita por el monitor `NS-MSG-Monitor`.
- **/(`serverId`).** Muy similar a la cola `/UDP`, pero con la diferencia que es leída por los diferentes servidores de WebSockets que hay en la infraestructura, o `NS-UA-WS`.

Debido a la alta carga, este servidor se despliega en modo cluster, con alta disponibilidad y mirroring de colas entre las diferentes instancias.

Base de datos, NO-SQL (MongoDB)

Un cluster de instancias de MongoDB es usado para mantener todos los datos persistentes que haya que guardar, como los registros de dispositivos y de aplicaciones, así como mensajes recibidos y los diferentes servidores de wakeup disponibles en las diferentes redes móviles.

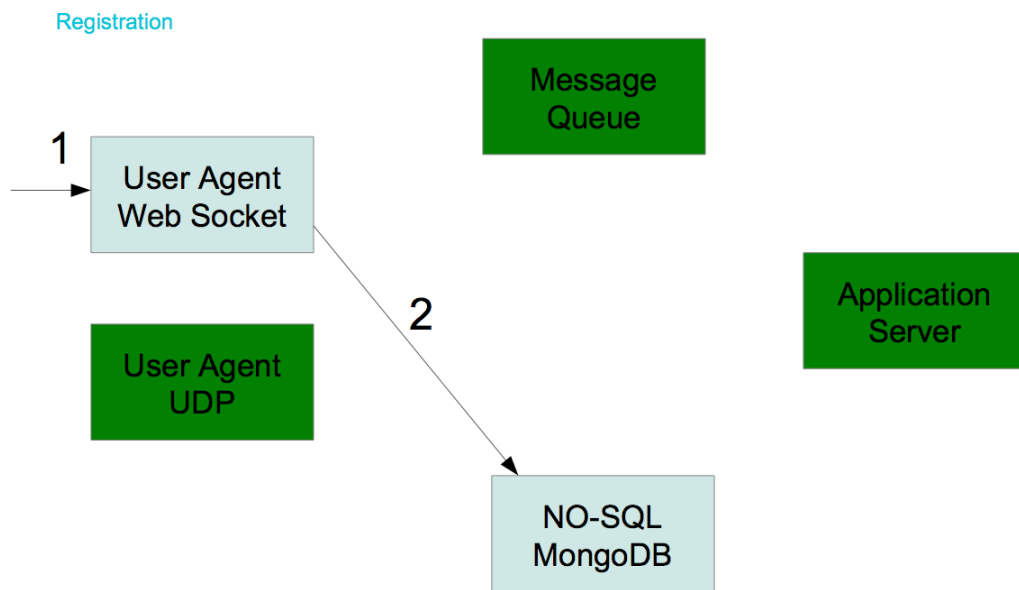
Flujos de trabajo

Así que una vez que están presentados los diferentes actores que contribuyen al funcionamiento correcto del servidor, se puede hablar de cómo se gestionan internamente los diferentes escenarios que pueden suceder. Así pues:

Registros de dispositivos/aplicaciones

Una vez hay disponible una conexión abierta, es flujo de registro es el siguiente:

Figura 4.2. Flujo de registro

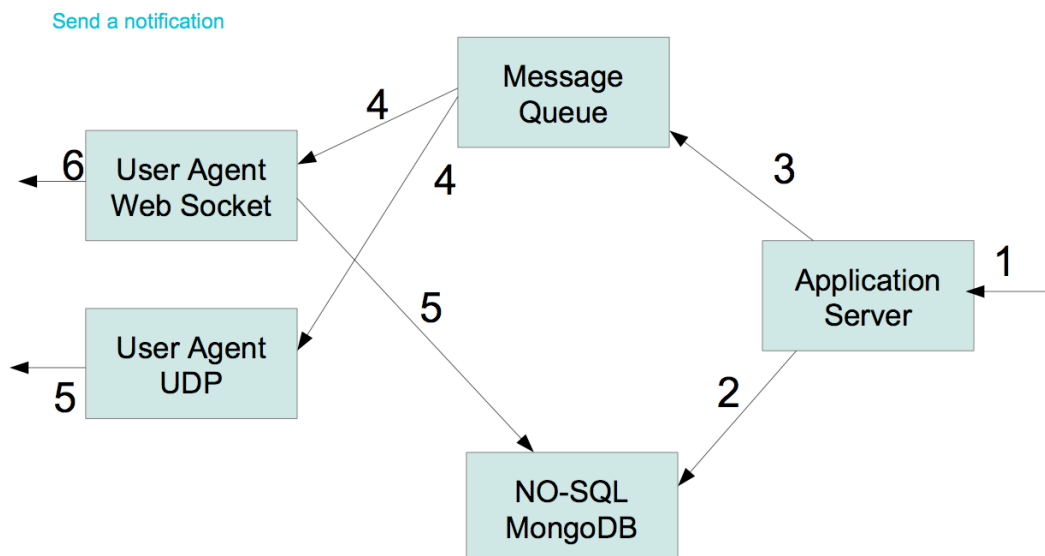


Como se puede observar, sólo actúan dos componentes del servidor de notificaciones, y es como sigue:

1. El registro es recibido por el servidor `NS-UA-WS`, que lo procesa.
2. Una vez procesado y aceptado, los datos necesarios, como la IP, puerto y demás información, como el token de usuario o de aplicación, es guardado en la base de datos para posterior uso.

Aceptación de notificación

Figura 4.3. Flujo de notificación



Según la figura anterior, los únicos puntos que se ponen en marcha son el 1, 2 y 3.

1. La notificación es recibida por el servidor `NS-AS`, que la procesa y verifica si es válida. Para ello, tiene que buscar si la URL a la que está destinada existe o no, haciendo una consulta a la base de datos. Además, también tiene que traer la llave pública de dicha URL para comprobar si la firma de la notificación es correcta o no.
2. Una vez verificada, se envía a la base de datos para que se guarde como método de backup, y para mantener la persistencia.
3. A la vez, se envía a la cola de mensajes, en concreto a `/newMessages`, indicando que ha entrado una nueva notificación en el sistema. Esta cola es leída por el servidor `NS-MSG-Monitor`.

Envío de notificación al dispositivo: WebSocket

Aquí ocurre el paso 4 de la figura anterior.

4. El `NS-MSG-Monitor` lee la cola `/newMessages`, preguntando a la base de datos la información necesaria sobre los nodos (o dispositivos) que tienen que recibir el mensaje. Dicha consulta devuelve una lista de colas a las que hay que enviar el mensaje.

Este paso comienza realmente cuando el mensaje ha sido enviado a la cola determinada que es leída por un servidor de WebSockets (`NS-UA-WS`) determinado.

5. Ahora que la notificación está en el servidor en concreto, se tiene que buscar en la base de datos quién es el dispositivo en concreto al que hay mandar.
6. Una vez recogida la información, se busca en memoria el conector que sea el que corresponda a los datos recibidos, y mediante el canal WebSocket abierto se manda la notificación.

Envío de notificación al dispositivo: UDP

4. El mensaje es recibido por la cola UDP, junto a una serie de datos que venían ya en la notificación, que le permite discernir a qué servidor de WakeUp tiene que enviar la petición para poder despertar al teléfono.
5. El servidor de WakeUp recibe una petición HTTP GET con dos parámetros: `ip` y `port`, que serán usados para enviar un datagrama UDP a dicha dirección para despertar al teléfono.

6. Una vez enviada dicho datagrama UDP, el teléfono se despierta y abre una conexión WebSocket con el servidor de notificaciones, enviando un mensaje de tipo `registerUA`
7. El servidor responderá con el mensaje especificado en capítulos anteriores, que incluirá un array con los diferentes mensajes que tienen que ser entregados.

Confirmación de recepción

La confirmación de recepción la realiza el servidor `NS-UA-WS` a través de una conexión de WebSocket que tenga abierta con algún cliente.

Para ello, se emplea el token de la conexión, ya que se guarda la primera vez que se hace el `registerUA`, y se espera un mensaje de tipo (`messageType`) `ack`. En este momento, el servidor pide a la base de datos que elimine de la colección de dispositivos, en el `uatoken` concreto, el mensaje con el id recién recibido. En términos de MongoDB, es hacer un `pull` de un elemento en una colección.

A partir de este momento, el mensaje ha desaparecido del servidor de notificaciones.

Capítulo 5. Tecnologías

Ahora, una vez se ha hablado de cómo es la arquitectura y la primera implementación, hay que comentar sobre las tecnologías usadas y por qué se han escogido para la parte del servidor, comparando con otras soluciones similares.

Así pues, una de las grandes decisiones que hay que tomar al iniciar un proyecto de software son las tecnologías y lenguajes de programación que se van a utilizar, incluyendo a las herramientas, editores, gestores de bugs y de incidencias, método de comentarios en el código y un sin fin de características que harán derivar el proyecto hacia el éxito o el fracaso.

Node.js

Node.js¹ es un nuevo paradigma de programación para aplicaciones web en la parte del servidor, basado en JavaScript, muy rápido y con una serie de características muy importantes, como se verá más adelante.

Node.js es una plataforma creada sobre el intérprete de JavaScript de Chrome para crear fácilmente aplicaciones de red rápidas y escalables. Node.js usa un sistema de eventos, un modelo de entrada-salida no bloqueante que hace que sea ligero y eficiente, perfecto para aplicaciones en tiempo real que requieran un uso intensivo de datos y que funcione a través de dispositivos distribuidos.

—Node.js

Esta nueva plataforma surgió en el año 2009 por Ryan Dahl, y es patrocinada por la empresa Joyent, y su creación se debió a la necesidad de crear un sistema que pudiera trabajar correctamente con sitios web que aceptaran peticiones push para el intercambio de datos, por lo que la idea de esta tecnología era especialmente interesante para el proyecto del servidor de notificaciones.

Está compuesta de diferentes paquetes independientes, que son:

- **Motor JavaScript V8 de Google Chrome.** Un intérprete creado ad-hoc para Google Chrome, que revolucionó a la competencia por ser el más rápido y usar unas ideas muy novedosas e inteligentes. Genera código máquina para varias plataformas, como 32-bits, 64-bits, ARM o MIPS antes de ejecutar el código, en vez de directamente interpretarlo como hacían otros motores JavaScript. Además, utiliza una serie de técnicas antes de ejecutar que permiten optimizar (en varios pasos)

¹<http://nodejs.org/>

y cachear resultados de partes del código que se usan mucho, como bucles o algunas funciones.

- **La capa de abstracción libUV.** Node.js se basó en la librería `libev` de abstracción de plataforma que estaba solamente disponible para Linux. Una vez que el desarrollo continuó y fue necesario utilizar otras plataformas, como Windows, se decidió crear `libuv`, que no es más que la evolución de `libev` para soportar otras plataformas y dar un API estándar para los desarrolladores de forma transparente al sistema que lo está ejecutando.

libUV contiene la mayoría de las características que luego definen a Node.js. Así pues, se implementan sockets TCP no bloqueantes, tuberías, el protocolo UDP, temporizadores, creación de procesos hijos e hilos, paso de mensajes vía IPC y sockets, varias funciones asíncronas como DNS o acceso a ficheros del disco...

Así pues, libuv obliga a un estilo de programación asíncrono y basado en eventos, y su trabajo principal es el de proveer un bucle de eventos y notificaciones basadas en callbacks para operaciones de entrada y salida y otras actividades.

- **Núcleo.** Escrito principalmente en JavaScript haciendo uso de las API e interfaces que provee V8 y otros módulos para un trabajo comunitario, es el código principal que después, mediante otras APIs se expone al desarrollador que va a utilizar Node.js.

Así pues, es posible encontrar módulos como el de criptografía² (`crypto`) escritos enteramente en C, haciendo uso de las extensiones de V8, mientras que por ejemplo el que provee la interfaz para el manejo de sistemas de ficheros³ (`fs`) está enteramente escrito en JavaScript, y su sintaxis es más fácil de entender y leer.

Para su diseño se tomó ideas prestadas de otros sistemas similares, como es Twisted para Python, libevent para C o EventMachine para Ruby. La mayoría de los programas en JavaScript se ejecutan en el contexto del navegador, puesto que siempre han sido los mayores exponentes y consumidores de este lenguaje, sin embargo, Node.js es ejecutado en la parte del servidor, e incluso puede ser usado como lenguaje de scripting (sustituyendo incluso a Bash o Python) y posee un REPL, una consola interactiva.

Usuarios. Actualmente Node.js, a pesar de ser una tecnología muy reciente y joven, está siendo usada por muchas grandes empresas de Internet para escalar sus sistemas. El hecho de que Yahoo! haya empezado a migrar algunos de sus sistemas desde PHP hacia Node.js,

²https://github.com/joyent/node/blob/master/src/node_crypto.cc

³<https://github.com/joyent/node/blob/master/lib/fs.js>

implementando diferentes frameworks y módulos mientras hacían la transición, hace pensar que el recorrido de esta plataforma es muy amplio y con muchas empresas grandes detrás interesadas en que siga siendo funcional, agregue más módulos y características y que sea líder en el mercado de los desarrollos de servicios de backend. Además, Microsoft, LinkedIn (para su aplicación móvil) o incluso Walmart también han decidido cambiar algunos de sus sistemas para hacer uso de Node.js.

Para qué se usa

Node.js se utiliza en la parte del backend como la programación lógica del servidor de notificaciones. Esto quiere decir que es el que controla todo lo que está ocurriendo en el servicio, expone las APIs de terceros, mantiene las conexiones abiertas, hace intercambio de datos, guarda en la base de datos, monitoriza los mensajes por si hay que borrarles y un sin fin más de funcionalidades.

Así pues, es el componente principal del servicio, el que está en continuo contacto con el mensaje enviado, desde la entrada, incluyendo la verificación, hasta la entrega a los diferentes terminales a los que iba dirigido. Es el maestro de ceremonias que hace que todo funcione correctamente.

Por qué se ha elegido

Hay multitud de tecnologías para programar servicios web, pasando por cientos de diferentes lenguajes de programación que proveen miles de frameworks diferentes entre los que elegir. Sin embargo, una tecnología robusta, sencilla, con mucha comunidad detrás para buscar cualquier duda y tener un montón de módulos ha hecho que el proyecto se decantara por Node.js por las siguientes razones:

- Es asíncrono y basado en eventos, lo que permite mucha más granularidad y no bloquear la ejecución. Además, no se tienen que lanzar hilos por cada petición (o varias), si no que está todo en un bucle de eventos al que se van añadiendo acciones.
- Es JavaScript, y está alineado con la tecnología del proyecto que lo engloba, que es FirefoxOS.
- Es muy rápido de programar y se pueden hacer pruebas de concepto de forma rápida, además, que sea JavaScript ayuda a detectar de forma más fácil los errores, y que el tiempo desde hacer un cambio hasta verlo reflejado sea mínimo, sin tener que esperar a grandes tiempos de compilación.
- Despliegue muy sencillo y con multitud de opciones, ya que permite elegir si se quiere correr detrás de un balanceador de carga (como

HAProxy), directamente conectado a la red y que sea Node.js el que responda...

Así pues, aunque hay otras elecciones posibles, como podría ser PHP, o incluso Python con pyev⁴, el que se creyó que tendría menos uso de CPU y memoria y tendría más posibilidades de ser "hackeado" por la comunidad, fue Node.js.

RabbitMQ

RabbitMQ es una plataforma de mensajería para aplicaciones, fácil de usar, que funciona en la mayoría de los sistemas operativos, de código abierto y con multitud de wrappers para diferentes lenguajes de programación, como Node.JS, C++, Java, Python...

Una plataforma de mensajería permite a las aplicaciones conectarse y escalar de forma sencilla y rápida. Las aplicaciones pueden conectarse con otras aplicaciones, o bien diferentes componentes pueden ser independientes y constituir entre todos una gran aplicación o incluso poder enviar o recibir datos. La mensajería es asíncrona, por lo que se desacoplan las aplicaciones y se basa en eventos (de recepción, de envío).

Algunas de las grandes características de RabbitMQ son:

- **Confiabilidad.** RabbitMQ tiene una serie de funcionalidades que hace que la entrega de los mensajes pueda ser confiable, mediante la persistencia de mensajes (el mensaje no se guarda sólo en RAM, si no también en disco, recuperable en caso de que haya un error en el sistema), aceptaciones de entrega (puede ser automático o manual, por ejemplo, se quiere decir que un mensaje está recibido cuando se ha procesado correctamente), confirmación de publicación (si ha habido fallos al enviar mensajes a las colas) y alta escalabilidad.
- **Enrutado flexible.** Los mensajes son enrutados a través de intercambios antes de llegar a las colas, y RabbitMQ provee varias técnicas de enrutado, como hacer round-robin, el típico publicador-suscriptor o incluso se pueden mezclar técnicas simples con complejas para poder realizar cualquier enrutado imaginable.
- **Clustering.** Se pueden conectar varias instancias de RabbitMQ en una red local, haciendo un único y simple servidor.
- **Federación.** Es similar al modo cluster pero con mayor tolerancia a fallos, usando redes que son más propensas a fallos, como en una

⁴<https://code.google.com/p/pyev/>

red WAN, por ejemplo, teniendo las instancias en diferentes centros de datos distribuidos por el mundo.

- **Colas con alta disponibilidad.** Las colas pueden crearse en modo espejo en varias máquinas en un mismo cluster, por lo que si una máquina falla, la cola estará replicada en otra, y el fallo será transparente para los programas que usen este sistema de mensajería. Además, automáticamente se comprobará y se intentará unir de nuevo los nodos caídos.
- **Cientes múltiples.** Como se ha comentado más arriba, RabbitMQ ofrece muchos clientes para lenguajes de programación populares, ya sea para desarrollos de backend, como C++, Node.js, Python, Ruby o incluso Erlang, como para clientes frontend, como JavaScript en el navegador.

Para qué se usa

El uso principal en el servidor de notificaciones es la de pasar mensajes entre los diferentes actores que están en ejecución, de forma desacoplada y funcionando de forma independiente, esperando a eventos que desencadenen una reacción dependiendo de qué contengan y qué haya que realizar según el mensaje.

El uso más simple que se da es el de algunas intancias creando contenido, que son los que publican y envían los mensajes a las colas, y otros que son los consumidores, los que reciben los mensajes y realizan con ellos diferentes tareas. Así pues, los suscriptores están apuntados a determinadas colas o tópicos, esperando a que lleguen mensajes que alguien les mete por el otro extremo. Es el caso de la cola `newMessages`, la cual es escrita por el servidor que provee la interfaz a los servidores de terceros, llamado Aplicacion Server (AS), metiendo, bien los mensajes completos (cuando su tiempo de vida sea 0), o bien el identificador que él ha creado y ha guardado en la base de datos, para que el consumidor lo lea, y entienda que tiene que entregar los mensajes dependiendo de algunos parámetros que se pasan en dicho mensaje en la cola.

Por qué se ha elegido

En un primer momento el software de mensajería utilizado era ActiveMQ⁵ ya que en el plan tecnológico se prefería el uso de esta tecnología, sin especificar las razones.

Sin embargo, dado el uso que se estaba dando en el servidor de notificaciones, en el cual no era requerido que hubiera una persistencia

⁵<http://activemq.apache.org/>

de datos, que era uno de los puntos flojos de ActiveMQ, se empezó a usar el protocolo STOMP⁶, que es el predeterminado de ActiveMQ.

El problema surgió cuando hubo que probar el despliegue, añadiendo pruebas de alta disponibilidad, tolerancia a fallos y durabilidad en el tiempo, momento en el que se descartó ActiveMQ debido a que era muy débil en cuando alguno de sus nodos fallaran, así que como el cluster no era como tal, si no que era solamente maestro-esclavo, lo que hacía que la escalabilidad fuera muy mala.

A partir de ese momento se buscaron otras soluciones de mensajería, viendo Apollo⁷, que es una nueva implementación de ActiveMQ basada en eventos y mucho más eficaz que su hermano mayor, ZeroMQ⁸, descartada por su gran curva de aprendizaje y su API a muy bajo nivel, y al final RabbitMQ, con el cual se hicieron más pruebas de carga, se decidió que sería la fundación final para los problemas de mensajería que se habían presentado, obligando a cambiar el protocolo de intercambio de datos, de STOMP a AMQP⁹, nativo de RabbitMQ y estandarizado e implementado por la mayoría de los sistemas de colas (por lo que en un futuro se podría cambiar el software de mensajería sin que diera problemas el servidor de notificaciones).



Nota

Hay que destacar que debido a los problemas que se encontraron en ActiveMQ el plan tecnológico está siendo revisado y es muy probable que en un futuro cercano se prefiera la opción de RabbitMQ en vez de ActiveMQ.

MongoDB

Es una base de datos NoSQL escalable, orientada al almacén de documentos, con alto rendimiento, de código libre y escrita en C++. Es rápida y provee una flexibilidad que no dan las bases de datos relacionales, sobre todo en la manera como se guardan los datos, ya que no están ancladas a esquemas fijos que no pueden cambiarse.

Creado por la empresa 10gen¹⁰ en octubre de 2007, utiliza JavaScript para las consultas (el API es accesible y se escribe como JavaScript), usando el motor de renderizado SpiderMonkey de Mozilla.

La forma de guardado de los datos es en formato BSON, o Binary JSON, por lo que es totalmente compatible con JavaScript y no hay que hacer un

⁶<http://stomp.github.com/>

⁷<http://activemq.apache.org/apollo/>

⁸<http://www.zeromq.org/>

⁹<http://www.amqp.org/>

¹⁰<http://www.10gen.com/>

post-procesado de la información para utilizarla de forma agradable en el backend, ya que es un tipo de dato nativo más para JavaScript, que es el lenguaje de programación de la parte de backend, con Node.js.

Así pues, MongoDB es una base de datos muy rápida, muy flexible y con cientos de parámetros de configuración que hace que funcione desde una sola instancia hasta cientos de máquinas trabajando en paralelo en modo cluster, haciendo réplicas o muchas más maneras diferentes.

Algunas de sus características son:

- **Esquemas libres.** Que se guardan en un formato JSON, por lo que son fácilmente parseables por la mayoría de los lenguajes de programación. Además, permite realizar cambios al vuelo a los documentos, agregando o eliminando contenido sin tener problemas con si el esquema se repite o no, porque no hay tal.
- **Soporte completo de índices.** Aunque se guarden los documentos sin esquema, se pueden crear índices en aquellos campos que vayan a estar siempre. Por defecto, el índice genérico se realiza sobre el campo `_id` que es como la clave primaria para MongoDB. Sin embargo, se puede indexar por cualquier campo, incluso haciéndoles combinados para buscar de forma más rápida.
- **Replicación y alta disponibilidad.** Una de las razones por las que MongoDB es muy usado es porque tiene multitud de opciones para replicar los datos entre diferentes servidores, ya sea en redes locales LAN o redes más grandes, WAN.
- **Sharding.** Más allá de poder replicar la base de datos entre diferentes instancias, se puede crear la opción de dividir la base de datos dependiendo de los datos de entrada, y dirigir las escrituras hacia un lado u otro. Esto es interesante por ejemplo si se necesita cumplir alguna ley que obligue que en un servicio mundial, los datos de los usuarios de algún determinado país esté en un centro de datos en concreto.
- **Modificaciones atómicas.** Algunas bases de datos NoSQL rompen las convenciones de las bases de datos SQL usadas desde hace tiempo, como la atomicidad (operaciones que o bien se ejecutan en su totalidad, o no se ejecutan), que podría llevar a pérdida de información e incongruencias. Así pues, MongoDB posee un motor que le permite realizar este tipo de opciones sin coste alguno.
- **Map-reduce.** El famoso algoritmo map-reduce, cada vez más utilizado para realizar operaciones sobre grandes conjuntos de datos, es un ciudadano de primer nivel en MongoDB, lo que permite realizar operaciones de agregación y procesamiento de datos de forma sencilla.

- **GridFS.** Permite guardar ficheros de cualquier tamaño en la base de datos, y recogerlos de forma sencilla, incluyendo replicación, sharding y todas las características que tiene MongoDB descritas más arriba.

Para qué se usa

Al ser una base de datos, se utiliza para guardar la información que envían los diferentes actores o instancias que están en ejecución en el servidor de notificaciones. Así por ejemplo, se guardan las notificaciones enviadas y correctas, los registros de los diferentes dispositivos con sus aplicaciones enlazadas y los posibles mensajes...

Por qué se ha elegido

Por proveer una flexibilidad mucho mayor que las bases de datos SQL tradicionales, y porque el formato de los mensajes, ya sean de los recibidos y los que se intercambian en las colas y de los datos entre las diferentes instancias tienen una característica en común: son JSON. Así pues, cualquier dato del sistema puede guardarse sin problemas en la base de datos, haciendo incluso de respaldo de los datos en caso de fallo del sistema.

Capítulo 6. Seguridad

La seguridad es uno de los puntos más importantes en cualquier sistema informático, y por supuesto no ha podido obviarse en este proyecto. Tener un sistema seguro es una tarea complicada, que requiere de mucho esfuerzo y dedicación, y sobre todo, de investigar cuáles pueden ser los fallos del sistema que pueden ser abusados. Aún así, cualquier sistema, por muy seguro que sea, siempre puede ser mejorado, y debe serlo, puesto que en el servidor de notificaciones se están moviendo mensajes, que pueden ir en claro, e identificando quién los manda y hacia qué terminales, por lo que puede ser un grave problema si se tiene algún fallo de seguridad.

Así pues, ha habido diferentes elementos del sistema en los que ha habido que poner más énfasis en la seguridad, bien porque eran susceptibles de mostrar información personal, bien porque se creía que podrían recibir ataques por fuerza bruta o de denegación de servicio o incluso porque son aquellos que están expuestos tanto a los terminales como a Internet, y siempre son un dulce muy apetitoso para personas con no buenas intenciones.

Tokens de dispositivo

El primer punto de entrada al sistema de notificaciones push es tener un token de dispositivo válido que permitirá que el UA (User Agent, explicado en el capítulo "API") se registre de forma correcta en el sistema.

Entonces, había que garantizar de algún modo que este token fuera tanto único como válido en nuestro sistema para poder controlar el abuso, y si los registros de dispositivos venían de donde donde se quería, esto es: de dispositivos con Firefox OS en un primer lugar.

A lo largo del tiempo, se han desarrollado muchos mecanismos de intercambio de información para realizar conexiones de usuario de forma controlada, o identificación de aplicaciones sin tener que dar los nombres de usuario y contraseñas a servicios de terceros, como OAuth¹. Sin embargo, en caso de la primera versión del servidor de notificaciones, se quería algo más sencillo, lo que se tradujo en los siguientes preceptos:

- **Verificación sencilla.** No debería requerir el uso de mucho cómputo para comprobar si el token de usuario es válido o no.
- **Rápido de generar.** Generación al vuelo, por lo que tendría que hacerse de forma rápida, ya que el cálculo era la respuesta a una

¹<http://oauth.net/2/>

petición HTTP y no se podía dejar dicha conexión HTTP abierta, por el uso de recursos en el sistema.

Posibilidad de cambiarse en un futuro. Si se decide a usar un nuevo sistema en el futuro, como pudiera ser OAuth, Mozilla Persona² (con lo que encajaría más aún en el sistema Firefox OS) o incluso que el token viniera preconfigurado desde fábrica no debería suponer un cambio drástico en la infraestructura.

Así pues, la decisión en la implementación de generar un token³, se quedó con el siguiente código:

Figura 6.1. Generación de token

```
function token() {}

token.prototype = {
  serialNumber: 1,

  // The TOKEN shall be unique
  get: function() {
    // SerialNumber + TimeStamp + NotificationServer_Id + CRC
    var rawToken = this.serialNumber++ + "#" + Date.now() +
      "#" + process.serverId + "_" + uuid.v1();

    // CRC
    rawToken += "@" + crypto.hashMD5(rawToken);

    // Encrypt token with AES
    return crypto.encryptAES(rawToken, cryptokey);
  },
};
```

Así, en la figura anterior, se puede ver cómo está generado el token de dispositivo:

1. Se usa un número de serie que aumenta en una unidad cada vez que hay una nueva petición.
2. Se añade el momento de generación actual, en milisegundos (según el estándar de JavaScript).
3. Se agrega el PID del proceso que lo ha generado.
4. Y un UUID (versión 1) para añadir aleatoriedad.
5. Se calcula el hash MD5 del `rawToken` inicial, que se lo añadimos al token inicial con un @ entre medias. Esto se usa como CRC.

²<https://support.mozilla.org/es/kb/que-es-y-como-funciona-browserid> . Mozilla Persona es un sistema de identificación similar a OpenID, con la diferencia de que Persona usa direcciones de correo en vez de URLs, lo que resulta más natural para la identificación

³Se puede ver el módulo que realiza esta generación y verificación en: https://github.com/telefonicaid/notification_server/blob/develop/src/common/token.js

6. Se devuelve el token generado cifrado con AES y la clave `cryptoKey` que está configurada en el fichero `config.js` que deben obligatoriamente compartir todas las instancias del servidor de notificaciones que generen tokens, pues es el principal medio para descifrar el token para realizar las operaciones y se usa para comprobar si es válido en el sistema.

Así pues, se puede observar que el token puede ser comprobado mediante la verificación con la función `verify(token)`, esto es, descifrando el token y comprobando su CRC. En primer lugar, si el descifrado no es correcto, el token recibido no es válido, por lo que se puede rechazar. Pero además, y en segundo lugar, es posible comprobar si el CRC (el hash MD5 calculado en la figura anterior) es válido, separando por el carácter `@` que habrá después de descifrar y comprobando si la segunda parte concuerda con la primera parte, habiéndole aplicado la misma función de hash.

Figura 6.2. Verificación de tokens

```
// Verify the given TOKEN
verify: function(token) {
  if(!token)
    return false;

  // Decrypt token
  var rawToken = crypto.decryptAES(token, cryptokey).split('@');

  // CRC Verification
  return (rawToken[1] == crypto.hashMD5(rawToken[0]));
}
```

Así pues, un posible ataque que podría intentarse es la suplantación de identidad, ya que sabiendo el token de dispositivo o `UAToken` es posible el registro de terceras partes como otro dispositivo. Sin embargo, esto no se considera un problema en el servidor de notificaciones y más en la parte cliente, aunque en todo caso, este token único tiene que estar gestionado por la seguridad que de el sistema operativo, o User Agent (UA) que implementa este protocolo.

Registro con clave pública-privada

Todas las aplicaciones que quieran recibir notificaciones tienen que registrarse con dos parámetros: el primero es un identificador único y el segundo una clave pública codificada en base64⁴, que es la encargada de la verificación de las firmas de cada notificación. Esto se realiza desde la propia aplicación web, con el API que hay entre la WA y el UA, mediante la función:

⁴RFC4648, página 5, punto 4, "Base 64 Encoding": <http://tools.ietf.org/html/rfc4648#page-5>

Figura 6.3. Registro de WA en el UA

```
DOMRequest navigator.mozPush.requestURL(DOMString watoken, DOMString pbk)
```

Así pues, es obligatorio que haya un par de claves públicas-privadas (soportando sólo `RSA-SHA256` en la primera versión) para poder registrarse y posteriormente enviar y recibir correctamente notificaciones.

Notificaciones firmadas

Como se ha comentado en la sección anterior, en el registro es obligatorio contar con una clave pública-privada basada en `RSA-SHA256`. Esto no se hace simplemente por añadir un atributo más al registro, si no por proteger a los dispositivos que están dentro de la red móvil. Así, esta idea se ha implementado para permitir notificaciones sólo desde lugares a los que se haya aceptado mandar notificaciones, esto es, aquellos registrados con una clave pública que puede verificar la firma de la notificación enviada con la clave privada que sólo tendrán los servidores de la aplicación (servidor `AS`).

El funcionamiento es el mismo que el sistema PGP con el correo firmado, ya que del contenido (atributo `message` de la notificación `JSON` enviada por `POST`) se crea una firma con la llave privada (que sólo la tendrá el servidor que envía la notificación, o `AS`), que se manda al servidor de notificaciones (a una URL determinada, que es capaz de saber la clave pública que verifica dicha firma) y comprueba si es correcta o no. Si lo es, deja pasar la notificación al sistema y, eventualmente, se entregará. Si no es válida, por la razón que sea, la notificación es rechazada y ni siquiera entra al sistema.

Figura 6.4. Verificación de firma (en `src/common/crypto.js`)

```
verifySignature: function(data, signature, publicKey) {  
  var algorithm = 'RSA-SHA256';  
  var verifier = crypto.createVerify(algorithm);  
  verifier.update(data);  
  return verifier.verify(publicKey, signature, 'hex');  
},
```

Figura 6.5. Comprobación de firma (en `src/ns_as/ns_as_server.js`)

```
//Get the PbK for the apptoken in the database
dataStore.getPbkApplication(apptoken, function(error, pbkbase64) {
  if (error) {
    return callback(errorcodesAS.BAD_MESSAGE_BAD_SIGNATURE);
  }
  var pbk = new Buffer(pbkbase64 || '', 'base64').toString('ascii');
  if (!crypto.verifySignature(normalizedNotification.message,
                             json.signature, pbk)) {
    log.debug('NS_AS:onNewPushMessage --> Rejected. Bad signature');
    return callback(errorcodesAS.BAD_MESSAGE_BAD_SIGNATURE);
  }
  ...
  // La firma es válida, continuar el flujo normal
  ...
}
```

Pero, ¿qué se consigue con esto? En primer lugar mantener acotados los servicios que pueden mandar notificaciones puesto que un desarrollador tiene control sobre la llave pública, puesto que es la que concuerda con su llave privada. Si pierde la llave privada, no podrá hacer nada, por lo que no podrá volver a mandar notificaciones. Pero además, si la recoge una persona malintencionada y decide empezar a mandar notificaciones causando una riada de mensajes, puede rápidamente actualizar su aplicación, poniendo su nueva clave pública en sus ficheros de configuración y la privada en sus servidores, y todos los dispositivos con esta nueva instalación se registrarán de nuevo (los dispositivos comprueban entre actualizaciones de la aplicación si la llave pública de la aplicación ha cambiado, y si es así, se registran de nuevo) por lo que no recibirán mensajes desde la llave robada, si no desde la nueva, ya que el sistema ni las dejará pasar.

Y en segundo lugar, se puede hacer un control de abuso llegado el caso. Por ejemplo, si un servicio no hace más que enviar mensajes de broadcast hacia el interior, se puede revocar todas los registros de aplicaciones que tienen dicha clave y se protegerá de forma fácil al usuario.

Verificación de notificaciones

En cuanto llega una notificación, o dicho de otra manera, cualquier mensaje mediante una petición `POST` a la interfaz que expone el servidor de notificaciones hacia el exterior (en este caso se habla del `API` hacia el `Application Server` o `AS`), hay que comprobar una serie de criterios para verificar que la notificación es correcta.

En primer lugar, se comprueba que el mensaje enviado es un mensaje `JSON` bien formado, con la función `JSON.parse(notification)` nativa de JavaScript. Si es así, se pasa a comprobar el resto de campos, de la forma:

- Que el atributo `messageType` de la notificación es igual a la cadena `notification`.
- Que esté correctamente firmado, como se explica en la sección anterior. Esto significa que tenga el atributo `signature` y además sea válido.
- Que tenga un identificador externo, esto es, que el atributo `id` no puede estar sin definir (`undefined` al parsear el `JSON`) o ser nulo, puesto que en un futuro se podría agregar la opción de sobrescribir notificaciones y poder preguntar por el estado de la entrega de una notificación a través de este identificador único para el servicio de terceros o `AS`.
- Que el atributo `message` e `id` no superen los 4KiB. Esto es debido a que la longitud máxima que se permite en el servidor es de 4KiB. Probablemente en el futuro se cambie la restricción de tamaño para la parte de `id` y se obligue a un tamaño más pequeño, ya que 4096 bytes de identificación suficiente para cualquier aplicación.

A partir de este momento, la notificación, si es válida, pasa al sistema que lo entregará a los dispositivos a los que iba destinada, eliminando algunos campos, como `signature` (para evitar ataques replay), y lo dejará de la siguiente manera:

Figura 6.6. Notificación normalizada

```
{
  "messageType": "notification",
  "id": 1234,
  "message": "Hola",
  "ttl": 2592000,
  "timestamp": "1356102444115",
  "priority": 1
}
```

Ataque DDoS: flooding y replay

Una manera muy fácil de colapsar sistemas informáticos expuestos a Internet es mediante ataques que inunden las comunicaciones y acaben con los recursos de forma rápida, mediante miles de ordenadores haciendo peticiones de forma continuada y de una manera distribuida, siendo imposible casi discernir cuáles bloquear porque la mayoría son peticiones válidas. Es por esto que hay que tener en cuenta muchas variables en el diseño de sistemas, a la vez que hay que minimizar el riesgo de lo que se deja abierto a terceras partes, y siempre monitorizar todo lo que está pasando para intentar identificar comportamientos extraños que pudieran entrañar un problema al sistema.

Una de las medidas tomadas es la normalización de las notificaciones (ver figura precedente), como se comenta en la sección anterior, eliminando la firma, ya que sería extremadamente fácil hacer ataques de repetición (replay) hacia el sistema, inundando todas las instancias con mensajes iguales pero totalmente válidos, al estar firmados. Así pues, aunque la notificación pueda ser enviada de nuevo, al no tener las aplicaciones o un servicio de terceros la firma del mensaje (que está compuesta gracias a la clave privada que sólo lo tiene el primer AS), el mensaje no se valida y es rechazado.

Comunicación cifrada vía SSL

Una de las maneras más extendidas por Internet para intercambiar información de forma segura es el uso de certificados reconocidos por alguna de las CA (Certificate Authority) de confianza. A pesar de que este sistema es bastante frágil⁵ y ha tenido numerosas⁶ debilidades⁷, sigue siendo el más usado y el que soportan la mayoría de los navegadores.

Es por eso, que todas las interfaces que expone el servidor de notificaciones hacia terceras partes, ya sean a los dispositivos (UA) o a los servidores de aplicación (AS) usan SSL (Secure Sockets Layer) para evitar que una posible persona malintencionada (normalmente conocido como Man-In-The-Middle o MITM) pueda interceptar la información y la descifre, con la obvia pérdida de privacidad.

Así pues, en la parte de los dispositivos se emplea dos métodos diferentes:

- **HTTPS.** La descarga del token inicial, en la primera versión del servidor, se realiza mediante un método GET a una dirección HTTP, por ejemplo `https://push.telefonica.com/token`. Esto permite que sólo los extremos de la conexión, en este caso el UA y el NS conozcan el contenido.
- **WSS (WebSockets Secure).** WebSockets también tiene una extensión segura para el intercambio de datos. Necesita un certificado SSL de la misma manera que una conexión HTTPS normal. De hecho, WebSockets se inicia mediante una petición HTTP, y después de hace un Upgrade para poder pasar el protocolo a WebSockets. Una vez negociada la conexión y actualizada a WebSockets, ya es posible transmitir información cifrada entre las dos partes. Por este canal se descargan las notificaciones, las llamadas de registro y confirmación de mensajes...

⁵BlackHat USA 2011: SSL And The Future Of Authenticity: <http://www.youtube.com/watch?v=Z7Wl2FW2TcA>

⁶Comodo creó certificados válidos para Google, Yahoo o Skype dirigidos a terceras partes: http://www.schneier.com/blog/archives/2011/03/comodo_group_is.html

⁷Detectado un certificado para dominios de Google emitido por una autoridad no competente: <http://arstechnica.com/security/2013/01/turkish-government-agency-spoofed-google-certificate-accidentally/>

UDP para notificaciones

Uno de los hechos diferenciadores del servidor de notificaciones, como se ha ido comentando durante todo el desarrollo de este proyecto, es que no es necesario mantener un socket o una conexión abierta entre el teléfono (UA) y el propio servidor para ser notificado de que hay información pendiente.

Así pues, se ha decidido que el paquete que se envía desde el servidor de notificaciones (NS) hacia el dispositivo (UA) sea del tipo UDP, por las siguientes cuestiones:

1. El paquete UDP, no llevará ningún tipo de contenido. Y si lo lleva, será descartado. Esto es debido a que no es necesario que haya corrección de errores, ni aceptación (también llamados `ACK`), no importa el orden (de hecho, no hay orden implícito), y tampoco es necesario un control de congestión.
2. Utilizar TCP hubiera hecho que la señalización en la red sea mucho mayor, ya que para hacer un simple "ping" para que el dispositivo (UA) se entere de que tiene una nueva notificación, sería necesario tener control de aceptación, retransmisión, temporizadores, orden, control de congestión... y en general demasiado cargado para algo que no requiere ningún `payload` para ser leído.



Nota

Probablemente en el futuro se añada algún mecanismo de información de cambio de puerto, por ejemplo: como es fácil mandar un paquete UDP a los diferentes dispositivos para que se despierten a un puerto por defecto (como es el caso de la primera versión, que está fijado en el 5000), estos automáticamente se conectan al servidor de notificaciones para descargar los mensajes, pero que, al ser un ataque, no habría. La idea es que si esto sucede, el dispositivo se re-registre de nuevo en el servicio indicando otro puerto, para que haya menos posibilidades de inundar el servidor.

Entrega de notificaciones a las aplicaciones

Un problema de seguridad que podría ocurrir es que haya dos aplicaciones instaladas en el dispositivo del usuario, y que ambas tengan permiso para recibir notificaciones push, pero una de ellas sea una malintencionada que se quiera hacer pasar por la original.

Así pues, el dilema viene cuando una aplicación es la verdadera y la otra es una que la intente suplantar. Esto quiere decir que la aplicación

web verdadera genera un token de aplicación (o `WAToken`) de una manera poco aleatoria y que es fácilmente adivinable. En este caso, si alguien sabe cuál es el algoritmo de generación de dicho token, podría replicarlo y registrarse de nuevo en el sistema como la aplicación "original". La duda en este caso surge de saber cuál es la aplicación original, cuál es la copia e incluso cuál se ha registrado en primer lugar. Así pues, la decisión es la de mandar la notificación sólo a la primera aplicación que se haya registrado con un token determinado, y rechazar el registro de sucesivas aplicaciones con igual identificador.



Nota

Sin embargo, esto podría cambiar en un futuro, pero como se comenta, hay demasiadas variables para poder elegir la aplicación real, puesto que no hay manera de discernir ahora mismo cuál es, o si una aplicación que intente suplantar a otra se haya registrado antes que la real.

Capítulo 7. Lecciones aprendidas

Durante el concepto inicial, desarrollo e implementación del servidor de push, se encontraron muchas variables que han hecho cambiar ligeramente el rumbo tomado por el proyecto. Sin embargo, la idea general seguía presente y los cambios que se fueron realizando eran debidos a cosas que se creía que eran más interesantes de tener, como mensajes de deregistro, posibilidad de encolar mensajes y mandar sólo un mensaje final, añadir `ACKs` específicos de aplicación y no usar los propios de TCP...

SQL contra NoSQL

Una de las primeras decisiones tomadas en el proyecto fue la de usar una base de datos no relacional, también llamadas No-SQL. En este caso se eligió MongoDB debido a que era la seleccionada por el plan tecnológico de la empresa. La decisión entre elegir una base de datos No-SQL frente a una SQL convencional (como puede ser MySQL u Oracle) fue simplemente la de no tener un esquema estricto para cada uno de los datos que se iban a introducir, ya que se creía que el esquema cambiaría en el futuro según se fuera iterando (y así sucedió de hecho), y permitiría un prototipado de la plataforma mucho más rápido.

Sin embargo, el cambiar radicalmente la mentalidad entre una solución SQL a algo No-SQL, muy relajado y sin un esquema fijo para todos los "elementos" de una misma "tabla" podría suponer un peligro, puesto que las nociones de SQL y muchas de las decisiones para estos sistemas se toman a partir de dichos conocimientos.

Así pues una de las primeras cosas fue "dejar de pensar de forma relacional", con todo lo que ello implica, como la pérdida de alguna de las propiedades ACID¹. E incluso la noción de tabla y de esquema.

Mientras que en las bases de datos SQL hay un montón de restricciones, los esquemas son muy fijos y hay que respetar todas las normas ACID, además de intentar normalizar la base de datos para poder cumplir las cuatro normas básicas, en NoSQL depende mucho de cuál sea la base de datos seleccionada el qué hacer, básicamente porque cada una tiene sus características y rompen las propiedades ACID de diferentes maneras.

Así pues, algunas de las características de la base de datos elegida, MongoDB, hacen que las consultas sean mucho más rápidas cuando esté todo en el mismo sitio, lo que quiere decir que es mejor duplicar datos a hacer varias consultas a la vez. De hecho, si hay que hacer varias consultas, es mejor que no sean complejas, y que se encargan

¹ACID: Atomicity, Consistency, Isolation y Durability

los diferentes backend de unir los datos, ya que así se liberan de tareas pesadas a los mismos gestores. Por ejemplo, hacer un `JOIN` en MongoDB es muy pesado, pero sin embargo es muy rápido si se hacen dos consultas por separado y las unimos en el código del programa.

Investigar nuevas tecnologías

Uno de los mayores retos del proyecto ha sido el uso de una tecnología nueva pero muy emergente y con una comunidad enorme detrás como es Node.js. Este lenguaje de programación permite hacer iteraciones muy rápidas, con muchos módulos para añadir características de una manera muy sencilla. Sin embargo, es JavaScript, lo que hace pensar a muchas personas que no es válido para desarrollo de backend, como sí lo podría ser PHP, Ruby o Python, expresamente diseñados para eso, ya que existe el concepto de que JavaScript sólo se ejecuta en el navegador de los usuarios, y su único cometido es el de manipular el DOM² de la página y poder hacer interacciones y cambios en tiempo real, incluso pudiendo hacer peticiones AJAX para agregar contenido de forma dinámica.

Sin embargo JavaScript es uno de los lenguajes más denostados y minusvalorados, ya que posee un montón de características que lo hacen especialmente rápido y, sobre todo, divertido de programar, amén de una velocidad de ejecución cada vez más rápida gracias a la competencia que se ha ido desarrollando los últimos años entre los principales desarrolladores de navegadores de Internet, que al final son los principales valedores de crear intérpretes de JavaScript más complejos, rápidos y eficientes.

La idea detrás de Node.js es que sea asíncrono y esté basado en un bucle de eventos, por lo que toda la programación hay que dirigirla hacia allí. Además, hay que tener en cuenta que JavaScript no tiene hilos, pero tampoco es necesario por la implementación del lenguaje (usando callbacks y eventos). Así pues, el cambio de mentalidad también es significativo, ya que es poco común en la programación estructurada tipo C o C++ el uso de eventos (aunque obviamente hay librerías para poder usarlos), mientras que en JavaScript, y en concreto en Node.js, lo son todo.

Así pues, muchas veces el código queda poco ordenado, debido a que la programación no es lineal, si no que puede ocurrir que una entrada empiece por la primera función, pero que tenga que llamar a la tercera que a su vez llama a la continuación de la primera para devolver la ejecución a la inicial. Sin embargo, pasada la diversión final con callbacks y eventos, Node.js ha destacado como un lenguaje muy potente, versátil y sobre todo, rápido y eficiente.

²DOM: Document Object Model

Tener tests unitarios y funcionales

Crear test y pruebas es algo que debería estar obligado en todos los proyectos, pero que, por unas cuestiones u otras, no lo suele estar, y se suele llevar a un plano secundario, muchas veces con la documentación. Sin embargo, es una de las piezas fundamentales en un sistema que esté en desarrollo. Dada la naturaleza de la informática, quizás cambios que se piensen que no afectan a otras partes, sí lo podrían hacer, y ahí es donde entran los tests.

Así pues, es fácil pensar en tener tests para prácticamente la totalidad del código y de las funciones que se van usando durante el programa. Simplemente es elegir la tecnología más fácil y sencilla para que dichos tests se ejecuten y que crearlos no sea algo que lleve más tiempo que el propio código que se quiere probar.

Así pues, en el caso del servidor de notificaciones, se ha empleado Node.js como lenguaje para crear los tests, tanto unitarios como funcionales, apoyándose en algunos casos de módulos específicamente diseñados para ello, como Vows, que permite ejecutar tests en paralelo, esperar a que unos terminen antes de que empiecen otros, y que se haga de forma muy sencilla y muy visual. Todo ello totalmente integrado en el fichero `Makefile` para ser ejecutado como `make tests` hace que la ejecución y comprobación de funcionalidades y de regresiones sea rápida y sencilla.

Priorización y enfoque en tareas

Uno de los principales enemigos de sacar adelante un proyecto es la falta de priorización de las tareas. Esto quiere decir que muchas veces lo más importante que hay que hacer para que un proyecto salga adelante no es lo mismo que lo que desarrolladores piensan. Quizas añadir una funcionalidad que va a usar el 5% de la gente y requiere el 20% del tiempo debido a su complejidad, podría priorizarse después de una que vaya a usar el 80% de la gente y que requiera ese 20% de tiempo.

Así pues, alguna de las funcionalidades que el servidor de notificaciones se han dejado para una futura versión 2 han sido fruto de estos acuerdos: cosas que son muy interesantes de tener, y que ayudarán de una gran manera a los desarrolladores que lo usen, pero que requieren demasiado esfuerzo para programarse que podría concentrarse en otras partes más fundamentales y que, en general, afectarán a más gente y harán incluso que esas nuevas features interesantes sean más fáciles y rápidas de implementar en un futuro.

Trabajar en abierto

Como se ha comentado anteriormente, la naturaleza de este proyecto era que fuera abierto y que pudiera ser usado por todo el mundo, a la vez que esperar nuevos colaboradores o personas que encontraran problemas y ayudaran en su arreglo hasta otras que les gustaría añadir nuevas funcionalidades y se encarguen de dichos nuevos escenarios.

Entonces, es sencillo entender que toda la tecnología, lenguajes de programación, compiladores, sistemas, programas sean abiertos. Pero es más, los sistemas donde se guarda el código fuente, los repositorios, también están abiertos, en la plataforma GitHub, lo que da mucha más visibilidad y facilidad a terceras personas de acercarse y ayudar.

Como se ha comentado anteriormente, los tests es una parte fundamental del desarrollo de una aplicación o servicio, así pues, aprovechando la excelente integración entre GitHub (el servicio donde se sube el código fuente) y el sistema de integración continua TravisCI, se decidió usarlos en cojunción para que cada vez que alguien suba un cambio en el código, automáticamente se genere una compilación y se lancen los tests, enviando correos y avisos a las personas que están configuradas en caso de que algo haya ido mal o se encuentren regresiones.

La importancia del hardware

Siempre se piensa que cuanto más capacidad y más rápido mejor, pero en el caso de sistemas reales también hay que tener en cuenta la ubicación de los servidores, la conectividad, la latencia y el presupuesto.

Lo más lógico es ver las diferentes necesidades de cada parte del sistema, y hacer pruebas de carga en diferentes escenarios: red congestionada, red disponible al 100%, red con fallos intermitentes, discos duros con poca velocidad, con mucha velocidad, poca RAM, mucha RAM...

Además, también es necesario probar de forma incesante todas las instancias. En el caso del servidor de notificaciones, se detectó que Node.js sólo permite que se ocupe 1.4GiB de memoria por cada proceso, por lo que sería necesario lanzar varios procesos a la vez del mismo servidor o usar otros sistemas para poder balancear la carga, como usando el módulo `cluster`³ disponible de forma nativa.

En el caso de MongoDB, las pruebas indicaron que el rendimiento de la base de datos mejoraba de una forma drástica si las instancias tenían poca memoria RAM pero un disco duro SSD, cuyo rendimiento es unas 10 veces superior a los discos duros magnéticos tradicionales.

³<http://nodejs.org/api/cluster.html>

Capítulo 8. Futuro: v2 y siguientes

Este capítulo está dedicado a todas aquellas ideas que se han quedado por el camino o que no son tan importantes para implementar en la primera versión del servidor de notificaciones, pero que podrían ser interesantes para futuras versiones o que se han pedido por otros proyectos (como pueden ser de mensajería instantánea).

Así pues, algunos de ellos tienen más prioridad que otros, algunos pueden tener más capacidad real de ser usados que otros, pero es posible que muchos de ellos lleguen a implementarse finalmente, y formen parte de la especificación estandarizada por el W3C.

Sobreescribir notificaciones

Una de las características más interesantes que permiten otros servidores de notificaciones es la posibilidad de sobreescribir notificaciones enviadas previamente en el tiempo pero que además no hayan sido entregadas, por lo que podría añadir más información, o bien cancelarse.

Ejemplo 8.1. Notificación sobre el desarrollo de un partido de fútbol

En el caso de una aplicación de resultados de fútbol, podría ser que cada notificación tuviera un estándar definido para indicar los goles, cambios de jugadores, tarjetas y demás incidencias. Podría ser que en un momento determinado, uno de los equipos marque un gol, y se requiera avisar a todos los usuarios para que la información se actualice. Sin embargo, se podría pensar en enviar un mensaje push con sólo la información determinada, pero, ¿y si alguien no ha recibido los anteriores? Requiere mucha más lógica por parte de la aplicación el recibir todos los mensajes, ordenarlos y pintarlos de una forma cronológica, cuando el orden de entrega no está garantizado.

Así entonces, la idea podría ser la de que el mensaje tuviera un identificador fijo para el partido, por ejemplo `valencia-madrid-liga-20130120` y que cada vez que hubiera un cambio que haya que enviar a los usuarios se añadiera más contenido, sobreescribiendo a los dispositivos que aún no hayan recibido la notificación con esta nueva (en vez de tener N mensajes, habría 1 mensaje con N eventos), y siempre tendrían lo último, independientemente del número de eventos que haya habido antes. En el caso de dispositivos que estén recibiendo notificaciones que ahora se sobreescriben, simplemente sería implementar la lógica de descartar la información anterior y pintar todo de nuevo, o bien hacer caso omiso de lo antiguo y utilizar sólo los nuevos eventos.

Recepción segura de token de usuario

Actualmente la recepción del token de dispositivo/usuario que se utiliza como primer paso para registrar a alguien en el servidor de notificaciones no es demasiado segura, ya que a todo el mundo que pide un token, se le asigna uno válido. Esto debería cambiar en un futuro e integrarse con otros sistemas de autenticación de usuarios, para dar tokens sólo a los usuarios autorizados.

Así pues, podría ser que en un futuro este token pudiera venir de forma pre-instalada en el propio dispositivo, que tenga que utilizarse una tecnología como OAuth2 para recoger el token o que incluso haya que identificarse como comprador del terminal y asignarlo a la cuenta de usuario para poder recibir notificaciones.

Sin embargo, este token depende mucho del entorno en el que se esté ejecutando el servidor de notificaciones, ya que este token no sale nunca del dispositivo, y podría asignarse de la mejor manera que el proveedor del servicio disponga, ya que la parte de cómo recibir el token no está estandarizada, ya que se supone que ha habido un intercambio seguro de información para recibirlo y el token que tiene el usuario/dispositivo es válido para todo el sistema.

Cambio en el tamaño máximo

En un primer momento, y como se comenta más atrás en el proyecto, el tamaño máximo de mensaje está restringido a simplemente 4KiB de información. Esta elección está hecha porque es la estándar de Google Cloud Messaging (GCM), y supone un tamaño aceptable para enviar datos a cualquier aplicación.

Sin embargo, esto podría cambiar en un futuro dependiendo del uso que se esté dando a la plataforma: si los desarrolladores prefieren mandar mensajes pequeños o grandes, si mandan el contenido cifrado, si existen otras características que aún no se están dando en las aplicaciones o en Internet...

Poder enviar números de versión

Este es una de las características fundacionales de Thialfi, el sistema definido por Google y que se usa en Google Chrome para la sincronización y los contactos. Sin embargo, el sistema actual del servidor de notificaciones tiene la opción de mandar números de versión, no de forma implícita, sino poniéndolo de forma explícita en el payload del mensaje enviado.

Así entonces en un futuro podría haber estas notificaciones:

```
{
  messageType: "notification",
  id: "<ID interno para el AS>",
  version: "<Número de versión>",
  signature: "<Firma>",
  ttl: "<Tiempo de vida>",
  timestamp: "<Desde EPOCH>",
  priority: "<Prioridad 1-4>"
}
```



Nota

Como recordatorio, el formato actual es:

```
{
  messageType: "notification",
  id: "<ID interno para el AS>",
  payload: "<Número de versión>",
  signature: "<Firma>",
  ttl: "<Tiempo de vida>",
  timestamp: "<Desde EPOCH>",
  priority: "<Prioridad 1-4>"
}
```

Como se puede comprobar viendo ambos mensajes, el primero es un subconjunto del segundo, ya que con el método actual, el servidor es capaz de enviar números de versión a las aplicaciones, construidas de tal manera que el mensaje que reciben no contiene un `payload` real, sino que es una versión que tiene que estar sincronizada con el servidor de la aplicación, y que se encargará de ir a realizar una segunda conexión para traer los datos reales, ya manejables por la aplicación.

Ejemplo 8.2. Uso de números de versión en aplicación del banco

Un uso real de números de versión podría ser la aplicación del banco, que quiere enviar que hay nuevos movimientos en la cuenta corriente, pero no desea enviar ningún dato personal a través de la red no gestionada por ellos, por lo que podría mandar un número de versión para que la propia aplicación se conecte con los servidores del banco, sin dejar ningún resto de información personal en el servidor de notificaciones, y descargue la información requerida para estar actualizado.

Implementar soporte para IPv6

El soporte de IPv6 debería ser una de las prioridades de la segunda versión del servidor, básicamente porque su uso es cada vez más extendido (más que por ser algo voluntario, por ser algo obligatorio), y puede tener características muy gratificantes para el servidor y de la manera que se envían las notificaciones.

La idea sobre IPv6 es la de aumentar el número de direcciones posibles para que cualquier dispositivo tenga una IP diferente a cualquier otro y por lo tanto sea único en la red, haciendo que pueda ser directamente accesible, sin tener que pasar por proxies intermedios que redirijan la petición hacia otras direcciones.

Esto es un gran avance que podría acabar directamente con la funcionalidad del servidor de notificaciones, ya que, al estar el dispositivo directamente en Internet, con una IP pública, y siendo accesible desde cualquier lado, no hay necesidad de un proxy intermedio (como es el servidor de notificaciones) o de mantener conexiones abiertas (como se hace en los otros sistemas) ya que los servicios que quieran enviar algo al teléfono lo podrían hacer directamente.

Sin embargo, hay un problema, y es que los proveedores de telefonía móvil pueden decidir varias cosas para que sus dispositivos no sean directamente accesibles por Internet:

- **Cortafuegos.** Para evitar que los dispositivos que están realmente conectados en la red se protejan ante posibles ataques, podrían implementarse sistemas de cortafuegos o firewall que permitirían cortar algunos tipos de conexiones que puedan ser perjudiciales para los dispositivos, como peticiones de PING, posibles ataques DDoS, que haría que la radio móvil esté siempre en el estado de máxima excitación y la batería de los terminales se gaste de forma rápida.
- **IP Privadas.** Aunque IPv6 provea miles de millones de diferentes direcciones y se pueda identificar a cada dispositivo por una IP diferente, las operadoras podrían decidir seguir creando redes privadas como ocurre ahora mismo en las redes IPv4, para proteger de la misma manera el acceso desde Internet hacia los teléfonos de forma incontrolada. Sin embargo, cumpliría lo mismo que haría el cortafuegos propuesto anteriormente, pero además, al estar en un segmento de red que no se puede llegar por ser privado, tendría el problema de que estos dispositivos no podrían ser contactados directamente, por lo que necesitarían un proxy intermedio para poder hacer las peticiones.

Así pues, aunque IPv6 podría pensarse que podría matar al servidor de notificaciones, como se ha comentado anteriormente, es muy probable que no lo haga, ya que en ambos casos se necesitarían reglas especiales en los firewall para permitir que sólo el rango de direcciones de funcionamiento del servidor de notificaciones pueda acceder a los dispositivos o bien tener una de las instancias dentro de la red móvil, como ocurre con el despliegue en IPv4, donde es necesario que haya algo dentro de la red móvil para poder hacer `PING` directamente a los terminales sin tener una conexión abierta y que sean capaces de responder a notificaciones.

Implementar prioridades

Esto podría ser una de las grandes nuevas funcionalidades en la segunda versión del servidor de notificaciones, ya que permite a los desarrolladores no molestar a los usuarios constantemente y esperar a que estén mirando el dispositivo para poder mostrar las notificaciones. GCM tiene un sistema similar, con la propiedad `delay_while_idle` en el JSON enviado, que permite que no se entregue hasta que el dispositivo esté de nuevo conectado a la red y el usuario realizando algún tipo de transferencia de datos.

Sin embargo, tener sólo dos estados quizás es demasiado poco para los desarrolladores, porque se podrían necesitar otros valores para diferentes prioridades, encolarlas y poderlas enviar de forma instantánea o bien esperar a que haya varias, no teniendo que ser todas del mismo origen, para enviar a la vez.

Ejemplo 8.3. Información sobre lugares cercanos y ofertas

Un uso interesante podría ser la posibilidad de tener una aplicación que informe a un servidor de cuál es la ubicación del dispositivo y quiera enviar ofertas de lugares cercanos a donde se encuentre, como bares, restaurantes o conciertos de música. Sin embargo, el recibir anuncios constantemente podría ser muy molesto para el usuario (incluso podría llegar a desinstalar la aplicación), por lo que se podría indicar al servidor cuál es la prioridad del mensaje: esto es, si se obliga a que se mande la notificación sea cual sea el estado del dispositivo, o que se pueda encolar y esperar hasta que el usuario esté de nuevo activo (por ejemplo, con la pantalla encendida o navegando por Internet) o se reciba una notificación de una prioridad superior, momento en que se enviarían todas.

PINGS de backup

Aunque la idea de este servidor de notificaciones es que no haga falta tener un canal abierto, constante, vacío, a la espera de hacer `PING` de vez en cuando para comprobar si hay datos, sí podría ser interesante en determinadas conexiones, como la de WebSockets.

La idea es tener este mecanismo de `PING` que se realice cada cierto tiempo (por definir, puede ser 10 minutos o media hora) y que aprovechando el `PONG` del servidor, pueda enviar la lista de registros de aplicaciones y los posibles mensajes, en el caso de que haya habido un fallo y no se haya podido entregar.

Sin embargo, esta es una de las tareas menos prioritarias porque se estaría volviendo a los sistemas anteriores que requerirían mantener un

canal abierto para preguntar de vez en cuando, el famoso polling, o bien abrir un canal para exclusivamente decir que se está conectado.

Soporte WAP push

El soporte UDP que se da en el servidor es nativo de la red: se usa toda la infraestructura montada para despertar al teléfono, que reciba el paquete, vea que es un `PING` de notificación, se conecte al servidor y descargue los datos.

Pero puede mejorarse con la introducción de WAP push, aunque sólo sería en casos muy excepcionales, porque ya se ha visto cómo funciona esta tecnología, abriendo el navegador del usuario para descargar los datos, algo que haría que se sintiera extraño, a menos que se pudiera hacer de fondo y se pudiera pasar a la aplicación en concreto.

La prioridad es baja por los problemas que tiene este sistema como se comenta en el párrafo anterior y porque es un servicio de pago para los desarrolladores o servicios que quieran mandar notificaciones push, además de que los usuarios tengan que dar su número de teléfono para poder hacerlo funcionar.

Control de abuso

Actualmente no hay ningún control de abuso para la primera versión del servidor. Se controlan los orígenes frente a una base de datos que se puede ir rellenando, pero no hay nada automatizado que lidie con todo el proceso.

Lo que se quiere para la segunda versión es la posibilidad de que esto sea automático, para que cuando el servidor detecte un alto número de conexiones, normalmente erróneas, sea suficientemente inteligente de ir añadiendo las IPs de origen a una lista negra y no permita que vuelva a conectarse en un futuro.

Sin embargo, no sólo sería para el control de quién se está conectando, si no incluso para aplicaciones que inundan el sistema con mensajes no válidos o que nunca serán entregados, o incluso que se envíen una y otra vez sin control.

Soporte a diferentes servidores de notificaciones

Uno de los grandes pilares de Internet ha sido siempre la descentralización de los servicios, pero la interoperabilidad entre ellos. Así por ejemplo ocurre con el correo electrónico, ya que las personas son

capaces de escribirse entre ellas aunque no sean del mismo proveedor (por ejemplo, esto no ocurre en los sistemas de mensajería instantánea de hoy día), ya que los protocolos están muy bien definidos y es capaz de enrutarse mensajes entre diferentes servicios, con la capacidad de que lleguen de unos a otros sin problemas.

Con esta misma idea ha nacido el servidor de notificaciones, queriendo que sea federado y que cualquier persona pueda instalarlo, incluso creando su propio servidor desde cero, mientras se respeten los estándares.

Así pues, la idea de que haya una red de servidores de notificaciones federado es algo lógico. Podría ser interesante que un servidor de Telefónica reciba alguna notificación que vaya a un servidor de Vodafone y sea suficientemente inteligente el enrutado para poder reencaminar la notificación sin molestar a la red interna de cada servidor.

Es por esto que quizás en versiones futuras se cambie la forma en que se da la URL de push (a la que el servidor de terceros tiene que enviar su notificación en formato JSON con un HTTP POST), de una URL normal, del estilo `http://push.telefonica.es/notify/abcdefghijklmnopqrstuvwxyz` a algo similar a una dirección de correo, y sea el propio sistema federado DNS el que sepa cómo reencaminar el paquete, del estilo de `abcdefghijklmnopqrstuvwxyz@push.telefonica.com` y que mediante DNS pueda resolver alguno de las entradas `SRV`¹ en el DNS donde se colocaría alguna información para enrutar el mensaje de forma correcta.

Control de presencia

Una de los principales argumentos de los desarrolladores de aplicaciones de mensajería instantánea es que requieren tener un control sobre si el usuario está conectado, desconectado o en estado ausente. Eso es sencillo para ellos si crean un propio socket contra sus servidores y mantienen la conexión y el estado ellos, sabiendo si se cae la conexión, usando APIs del sistema saber si la pantalla está encendida y avisar del estado...

Sin embargo, el principal problema es que los usuarios no suelen tener solo una aplicación de este tipo, si no que suelen tener varias, muchas veces dependiendo del círculo de amigos, o de donde tengan más contactos. Así pues, podría ocurrir que ese socket se multiplicara por dos, tres, o incluso 4, según el número de aplicaciones que lo requieran.

Es por esto que una de las ideas es la de proveer un API a los servicios para comprobar si alguien está conectado/desconectado o en estado

¹<http://www.ietf.org/rfc/rfc2782.txt>

ausente. Esto se podría hacer de una forma muy sencilla tal y como está montado el servidor de aplicaciones, ya que sería hacer una petición a la URL que se quiere saber si se puede entregar, se miraría el estado del nodo o de los nodos que corresponden a dicha URL y se devolvería cuál es el estado actual: conectado, desconectado, ausente o indefinido (cuando se esté bajo UDP y no haya una conexión abierta)

Esto trasladaría la carga de la red desde los terminales y la red de telefonía al servidor de notificaciones y a Internet, porque no tienen que estar respondiendo los terminales continuamente, ni manteniendo canales abiertos, si no que es el propio servidor el que sabe si el dispositivo puede ser alcanzable o no.

Control de entrega de notificaciones

Muchas de las aplicaciones de mensajería instantánea requieren saber si el mensaje enviado ha sido entregado o no. Podría ser para mostrar un pequeño icono al usuario para que sepa que el otro lado ha recibido el mensaje, o bien porque se quiere llevar el control de si se entrega el mensaje correctamente o si el dispositivo está conectado.

Así pues, este sistema proveería un API para poder preguntar sobre el estado de entrega de un mensaje, no sólomente si se ha entregado o no, si no quizás en qué momento, y si no se ha entregado, por qué razón: 2aplicación desinstalada, el dispositivo está desconectado...

Ejemplo 8.4. Enviar mensajes sólo si los anteriores han llegado

Un escenario típico para algunos servicios es la de que si no se han entregado los mensajes anteriores, no tiene sentido enviar mensajes posteriores. Esto podría ser válido en aplicaciones que manden un solo `PING` al usuario. Probablemente no sea interesante mandar varias notificaciones, si no saber si la anterior se ha entregado, y si es así, mandar un nuevo mensaje. En el caso de que no se haya entregado, en algún momento se entregará y el mensaje siguiente sería el mismo, por lo que no se envía, no ocupa espacio en el servidor de notificaciones, y reduce la carga en la red móvil, al sólo tener que enviar un único mensaje.

Capítulo 9. Conclusiones

Durante el desarrollo de este proyecto, se ha explicado cuál es el problema principal de las actuales soluciones de notificaciones push que están disponibles en los diferentes sistemas operativos móviles, viendo que la sobrecarga en las redes por usar conexiones permanentes es muy perjudicial, así como debido al alto número de mensajes de señalización que son necesarios para que funcione el protocolo GSM correctamente. Esto repercute negativamente en la salud de las redes, debido a que se congestionan con facilidad, y al final se están transmitiendo más datos de control que reales, lo que puede llevar a que no haya ningún dato usable por las aplicaciones en estas comunicaciones.

Así pues, la idea es la de usar todos los beneficios que proporcionan las redes móviles y los dispositivos que hacen uso de ellas, para realizar un sistema push más eficiente y amigable con la red, pensado para que incluso cualquier sistema operativo moderno pueda implementarlo y sea usable e interoperable.

Para ello, se han discutido los diferentes sistemas actuales de push en las plataformas móviles más conocidas, como son Android y iOS, viendo cuáles son sus puntos débiles y sus virtudes, para evitar tener los mismos problemas y a la vez coger lo más positivo para enriquecer el sistema aquí propuesto. Se ha visto que los sistemas son muy similares en su uso y concepción, y al final siempre mantienen una conexión abierta, con PINGs periódicos para hacer que no se caiga dicha conexión y poder mantener siempre un canal abierto entre los dos extremos.

Así pues, el aprovechar todas las posibilidades que tienen las redes móviles, conocer su funcionamiento y sus peculiaridades, ha ayudado a crear esta solución que ahorra en costes para las operadoras, puesto que tienen menos carga en su red (medida en mensajes de señalización y en conexiones abiertas), y repercute en los usuarios en un menor uso de batería y de datos, lo que influye en una mejor calidad del servicio y en un ahorro económico.

Además, el trabajo se ha realizado con vistas a la estandarización por parte del W3C, por lo que el ánimo es que esté disponible en todos los navegadores y que, dado su protocolo simple, también pueda implementarse en sistemas operativos, principalmente móviles, para tener ventaja sobre el resto, pero también en sistemas de escritorio, y permitir a los desarrolladores tener una única API con la que trabajar, realizando desarrollos más rápidos, eficientes y multiplataforma.

Un dato muy importante es la capacidad de que las notificaciones planteadas pueden transmitir datos, lo que influye en que las aplicaciones que las reciban puedan usarlas directamente, sin tener que hacer una

segunda conexión a su servidor, lo que repercutiría muy negativamente en las redes y en la batería, eliminando cualquier mejora planteada por las soluciones al no mantener canales abiertos con mensajes de `keep-alives`.

La primera implementación se ha realizado con éxito en los dispositivos con FirefoxOS, pudiendo realizar todo lo que se puede hacer con otros sistemas de push, como el de Android o el de iOS, incluso despertando a aplicaciones que están cerradas y tienen que ser abiertas para que procesen las notificaciones.

Como se ha demostrado, este servidor es tanto continuista como innovador, ya que utiliza algunas de las ideas de otros sistemas, como cuánto tamaño se permite en cada mensaje y el tipo de conexión, pero sin embargo, se añade la posibilidad de que el terminal no tenga conexiones abiertas continuamente y pueda ser despertado utilizando las diferentes ventajas que nos permite la red móvil.

Es por eso que puede suponer un gran cambio en el uso de los servicios de push en los diferentes sistemas operativos, debido a que es muy fácil de implementar y ayuda enormemente a las operadoras a que no tengan sus redes colapsadas y que los usuarios no gasten tantos datos en paquetes de control y batería por tener que mantener conexiones abiertas sin utilidad clara.

Por tanto, el futuro de esta tecnología es prometedor, puesto que ayuda a todas las partes implicadas en las comunicaciones móviles con intereses en las notificaciones push, a tener un mejor sistema que elimina mensajes redundantes, haciendo que las redes se comporten de manera más óptima, utilizando un sistema estándar para todos los sistemas, y que repercutirá muy positivamente en los usuarios, puesto que verán que la batería de sus dispositivos portátiles aumentará en tiempo y los datos descargados serán menores.

Glosario de términos

M

MongoDB	Base de datos no relacional que guarda sus datos en un formato BSON, similar a JSON, y cuyo lenguaje de consulta es JavaScript, teniendo funciones de reducción, mapeo...
---------	---

N

Notification Server	Servidor de notificaciones. Es el conjunto general en la parte centralizada, está compuesto de diferentes partes.
Node.js	Plataforma de desarrollo en la parte servidora escrita principalmente en JavaScript. Permite una programación rápida, basada en eventos y asíncrona.

R

RabbitMQ	Sistema de mensajería escrito en Erlang, muy rápido, con muchas configuraciones para las colas y preferencias para alta disponibilidad y alta carga.
Notification Server	Servidor de notificaciones. Es el conjunto general en la parte centralizada, está compuesto de diferentes partes.

U

User Agent	Agente de usuario. Es el programa que ejecuta una instancia web. En este caso, podría ser el propio sistema operativo Firefox OS o el navegador Firefox.
------------	--

W

WebApp	Aplicación web. Programa instalado en el dispositivo, basado en tecnologías web.
--------	--

Bibliografía

PDFs

[MAR12] *An Introduction to libuv*. Nikhil Marathe. Copyright © 2012 Nikhil Marathe. Libro online [<http://nikhilm.github.com/uvbook/An%20Introduction%20to%20libuv.pdf>]

Libros

[MAR12] *"An Introduction to libuv"*. Nikhil Marathe. Copyright © 2012 Nikhil Marathe. Libro online [<http://nikhilm.github.com/uvbook/An%20Introduction%20to%20libuv.pdf>]

Artículos

[MAR12] Nikhil Marathe. Copyright © 2012 Nikhil Marathe. *An Introduction to libuv*. Libro online [<http://nikhilm.github.com/uvbook/An%20Introduction%20to%20libuv.pdf>]